

Machine Learning for Administrative Health Records: A Systematic Review of Techniques and Applications

Adrian Caruana^{a,*,1}, Madhushi Bandara^{a,1}, Katarzyna Musial^b, Daniel Catchpole^{a,c} and Paul J. Kennedy^{a,d}

^aAustralian Artificial Intelligence Institute, Faculty of Engineering and IT, University of Technology Sydney, Australia

^bComplex Adaptive Systems Lab, Data Science Institute, Faculty of Engineering and IT, University of Technology Sydney, Australia

^cBiospecimen Research Services, The Children's Cancer Research Unit, The Children's Hospital at Westmead, Australia

^dJoint Research Centre in AI for Health and Wellness, University of Technology Sydney, Australia and Ontario Tech University, Canada

ARTICLE INFO

Keywords:

Machine Learning
Administrative Health Record
Health Informatics
Systematic Review
Pattern Mining
Population Health

ABSTRACT

Machine learning provides many powerful and effective techniques for analysing heterogeneous electronic health records (EHR). Administrative Health Records (AHR) are a subset of EHR collected for administrative purposes, and the use of machine learning on AHRs is a growing subfield of EHR analytics. Existing reviews of EHR analytics emphasise that the data-modality of the EHR limits the breadth of suitable machine learning techniques, and pursuable healthcare applications. Despite emphasising the importance of data modality, the literature fails to analyse which techniques and applications are relevant to AHRs. AHRs contain uniquely well-structured, categorically encoded records which are distinct from other data-modalities captured by EHRs, and they can provide valuable information pertaining to how patients interact with the healthcare system.

This paper systematically reviews AHR-based research, analysing 70 relevant studies and spanning multiple databases. We identify and analyse which machine learning techniques are applied to AHRs and which health informatics applications are pursued in AHR-based research. We also analyse how these techniques are applied in pursuit of each application, and identify the limitations of these approaches. We find that while AHR-based studies are disconnected from each other, the use of AHRs in health informatics research is substantial and accelerating. Our synthesis of these studies highlights the utility of AHRs for pursuing increasingly complex and diverse research objectives despite a number of pervading data- and technique-based limitations. Finally, through our findings, we propose a set of future research directions that can enhance the utility of AHR data and machine learning techniques for health informatics research.

1. Introduction

Machine learning research is gaining huge popularity in the health informatics domain, with researchers and practitioners alike interested in utilising data-driven technologies to assist in day-to-day clinical activities. Electronic Health Records (EHR) are the clinical data repositories underpinning this research. EHRs comprise heterogeneous data modalities that can pertain to any data that is recorded in a healthcare setting. The volume of EHRs is growing rapidly, with worldwide clinical data estimated have exceeded 2 yottabytes by 2020 with an annual growth rate of 48% [1]. Researchers, practitioners, and policymakers often use novel data analytics and machine learning techniques to understand healthcare processes, conduct clinical audits against national standards, plan future services and resource allocation, and for clinical governance [2].

Due to the heterogeneous nature of EHRs, the analytics techniques and healthcare applications pursued by EHR-based research depends strongly on the specific EHR data that is being analysed. For example, medical images, clinical notes, and time-series data from sensors are each common types of EHRs, yet require vastly different analytics techniques, and are useful in different aspects of healthcare research. This diversity has been established in reviews of

EHR data mining applications [3, 4], and other surveys that review specific methods for analysing EHRs [5, 6, 7, 8, 9, 10, 11, 12].

Administrative Health Records (AHR) refer specifically to the subset of EHRs that are collected for administrative purposes and comprise valuable information relating to how patients interact with the healthcare system. AHRs also have many unique characteristics that are distinct from other data-modalities captured by EHRs. AHRs contain extensive patient cohort information in the form of well-structured, categorically encoded records. In contrast, other EHR data modalities often contain unstructured or subjective records, such as clinical notes, lab reports, and medical images.

There are three significant gaps in existing EHR-based research. Firstly, despite the importance and unique characteristics of AHRs, existing EHR-based research often does not distinguish AHRs from other EHR data-modalities, nor does the existing research examine which data-driven methods are most applicable AHRs. Secondly, the standardised, population-level context of AHRs enables insight into many unique health informatics applications that are distinct from those identified in existing surveys. Finally, studies that analyse AHRs are often disconnected from each other, making it difficult for researchers to consolidate information from other relevant studies. These gaps hinder the development of novel analytics techniques for analysing AHRs and for

ORCID(s): 0000-0001-7283-0220 (A. Caruana)

¹Equal contribution.

pursuing important population-based health informatics research questions.

The objective of this survey is to identify state-of-the-art machine learning research applied to AHRs, with a focus on studies that seek to identify population-level health service patterns. This study systematically analyses 70 papers selected across multiple databases to synthesise and answer four key research questions. Using the results of this review, we subsequently propose a set of future research directions that can enhance the utility of AHR data for population-level health service pattern discovery. Finally, this review seeks to add clarity to the burgeoning research area of AHR-based analytics.

The review is organised as follows: Section 2 provides background on the current research landscape of mining AHRs for health service patterns, Section 3 details the systematic literature survey methodology we followed, Section 4 presents the answers to the research questions, Section 5 proposes some open research questions, Section 6 discusses the findings of this review, and Section 7 concludes the review.

2. Preliminaries

2.1. Background

We examined several recent survey and review papers that explore the intersection of machine learning applications and AHR data to determine the scope and extent of this review. Among them, Yadav et al. [3] and Chen et al. [4] conducted comprehensive surveys on a wide range of data mining techniques used for EHRs. Other literature reviews explored the application of specific techniques to EHRs, including process mining [5, 6, 7], network analysis [10], and deep learning [11, 12].

A fundamental gap in EHR-analytics literature is that AHR data is seldom distinguished from the many non-administrative data modalities comprising EHRs. Consequently, it is not clear which of the machine learning techniques identified in EHR-analytics reviews are suitable for analysing AHRs, or suitable for pursuing research objectives that pertain to AHRs. Since AHRs are used for many research objectives, (Section 4.3 details these objectives) this is a non-trivial gap in EHR-analytics literature.

To our knowledge, there are no existing studies that explore the broad landscape of methods suitable for analysing AHRs. This systematic review seeks to address this gap by choosing to focus on AHRs as the data source. The remainder of Section 2 introduces some useful definitions to contextualise the analysis of AHRs, and outlines the scope of analytical methods that will be explored in this systematic review.

2.2. Definitions

2.2.1. Administrative Health Records

According to Google Scholar, the term ‘electronic health records’ has appeared in over one million articles, while the term ‘administrative health records’ has only appeared in approximately one thousand articles. The use of AHRs

in literature is likely much higher than the search term ‘administrative health records’ suggests since many studies do not use this term even when the EHRs of concern are administrative. This ambiguity is problematic, since EHRs do not exclusively refer to administrative data, but encompass many other non-administrative modalities of healthcare data.

Cadarette and Wong [13] provide an introduction to administrative healthcare data from a pharmacology perspective, describing five common administrative databases and variables from Ontario. We use their descriptions, as well as examples of AHRs from Australia [14] and the United Kingdom [15] to inform the following definition and characterisation of AHRs.

Administrative health records are any data that is generated during interactions between various entities within a healthcare system, including but not limited to patients, physicians, hospitals, pharmacies, or government bodies. AHRs have the following five key characteristics:

- AHRs principally record attributes that are *discrete* or *categorical* in nature. Common examples include patient attributes, diagnosis or drug codes, hospital admissions, insurance claims, or death records.
- The attributes recorded in AHRs are typically *temporal*, meaning that they are recordings of events that occurred at a given time.
- AHRs contain linkage variables, such as patient ID, physician ID, or hospital ID, which facilitate the linkage of multiple AHRs.
- AHRs are typically *multivariate*, meaning that they can contain multiple attributes. For example, patient records may contain both diagnosis and treatment data.
- AHR attributes may also be *hierarchical*. This is common of attributes that are well-defined, highly structured, or when the set of values for the attribute is large.

2.2.2. Population Health Service Pattern Mining

AHR-based research often describes the process of learning healthcare patterns from AHRs in many different ways. In some cases, authors contextualise their work within a broad research context, such as process mining or healthcare informatics. However, most studies simply emphasise the techniques used or the objectives sought. These studies lack a shared vocabulary despite the common goal of learning about healthcare patterns from AHRs.

To clarify this, we use the term *health service pattern mining* (HSPM) to describe this shared goal. This term is a combination of both *health informatics* and *data mining*, and describes the process of analysing patient-level event data to gain insights into the operational processes within healthcare. Section 4.3 details some of the applications of HSPM, and Section 4.4 analyses how HSPM is done in practice.

We also refer to Kindig and Stoddart [16] for a definition of *population health*. Population health studies often use AHRs to pursue research questions concerning healthcare policies and interventions. This may be achieved by pursuing an intermediate research objective that is measurable, such as analysing health outcomes or patterns of care. An improved understanding of population health can be achieved by applying HSPM techniques to these measurable intermediate objectives.

2.3. Technique Scope

The scope of techniques explored in this review is informed by the definitions outlined in Section 2.2. In many cases, distinct analytical techniques share similarities in the way they represent, model, or evaluate data. However, in this review we distinguish between techniques specifically by their *modelling* approach. Consequently, the modelling techniques we explore may share commonalities. For example, this review specifies Bayesian networks and Markov models as distinct techniques, despite each using graphs to represent data. The remainder of this section uses this schema to outline and justify which specific analytics techniques fall within the scope of the review.

2.3.1. Machine Learning

The majority of the identified studies employ techniques that fall under the definition of *machine learning*. These methods include clustering, neural networks, regression analysis, Markov models, topic modelling, ensemble learning, Bayesian networks, association rule mining, sequence mining, temporal signature mining, and dimensionality reduction. Therefore, this systematic review focusses on machine learning as the principal method for analysing AHRs.

2.3.2. Process Mining

From our research, we observed that process mining was a widely explored technique for analysing EHRs, with multiple survey studies detailing its applications in healthcare process analysis [5, 6, 7, 8, 9]. Despite widespread adoption, many studies have identified significant limitations with process mining for EHR analysis.

Guzzo et al. [8] describe two significant limitations with this approach. Firstly, process mining struggles to accurately identify processes in EHRs since the EHR collection methods are not process-aware. Secondly, process discovery algorithms can struggle to generate interpretable process models from EHRs since they can vary significantly across patients. These studies suggest that AHR data are both less structured and more complex and diverse than the common business processes [17]. Most process mining algorithms are based on the assumption that the underlying processes happen in a structured fashion. Rebugue and Ferreira [17] suggests that this is a weak assumption in the context of human-centric clinical pathways; AHRs yield unstructured patterns and are often governed by decisions of individual patients and health practitioners, not by strict business processes. Furthermore, process mining is a technique more relevant to the domain of business process analysis and is based on rule-based,

black-box tools rather than machine learning techniques. Therefore, we excluded process mining based studies from the survey.

2.3.3. Network Analysis

Unlike other techniques examined in this review, network science is not a machine learning technique, but rather a means of representing discrete data and a suite of metrics and algorithms that operate on such a representation [18]. In our research, we identified many studies that utilise techniques from network science in pursuit of HSPM. This is because network science-based models can naturally capture healthcare events interconnected by varied and diverse relationships.

In addition to the representational benefits offered by networks, metrics and algorithms from network science can also be used to understand the characteristics of AHRs. While these algorithms are often not powerful enough to capture complex health service patterns embedded within AHRs, they can often capture simple, informative, and interpretable patterns (see the *Usage* section in Table 4 and Section 4.2.3).

For these reasons, we saw it prudent to include *network analysis* as a separate technique that includes the set of metrics and algorithms which operate on networks. This technique is not to be confused with other, more complex techniques that use graph-based representations, such as Markov-based techniques, Bayesian networks, or graph neural networks.

3. Research Method

Our study has been designed to answer the four research questions outlined in Section 3.1. These questions were used to inform the process of systematically reviewing relevant literature, including the bibliographic search process, inclusion and exclusion criteria, and the structure of the paper.

3.1. Research Questions

1. What are the machine learning techniques utilised in analysing health service patterns?
2. What are the applications of health service patterns in identified studies?
3. What are the strengths and weaknesses of specific machine learning techniques for analysing health service patterns?
4. What are the limitations of existing machine learning techniques for health service pattern discovery?

3.2. Bibliographic Search Process

To understand the research landscape, design keywords used for literature search, as well as synthesise the research questions, we followed an informal keyword-based search on *Google Scholar*, *IEEE Xplore*, and *Springer Link*. This step helped the authors understand current state-of-the-art and identify existing literature reviews. Through this process, we identified 42 studies. By analysing the frequent words contained in titles, keywords, and abstracts of these studies, we

came up with the following search string for the systematic literature search:

```
('data mining' OR 'machine learning' OR
'pattern mining' OR 'pattern recognition') AND
(patient OR treatment OR clinical OR healthcare) AND
(flow OR journey OR process OR pathway) AND
('electronic health record' OR
'administrative health record' OR
'Administrative healthcare record' OR
'electronic medical record')
```

Even though AHRs represent a special subset of EHR data, in many cases existing literature does not make this distinction clear. Therefore, we include 'electronic health record' and 'electronic medical record' terms in the systematic literature search as to ensure that all relevant articles are included. Studies that explore EHR data that is not administrative in nature are filtered out later in the search process.

We followed the process proposed by Petersen et al. [19] and Harris et al. [20] and conducted the initial evidence search on five databases (*IEEE Xplore*, *ACM Digital Library*, *ScienceDirect*, *PubMed*, *Springer Link*) for published research whose title or abstract meets the search string criteria outlined above. Publications of the *Artificial Intelligence in Medicine* was also searched for matching studies and are reported under *ScienceDirect* results. These primary sources were selected since they are prominent publications for computer science and machine learning as well as the healthcare domain. Finally, *dblp* was also used to search for any literature that may have been missed by the six main databases and also to ensure the latest studies not indexed in initial databases are also included as our evidence. Findings were further extended through snowballing approach proposed by Wohlin [21].

3.3. Filtering of Bibliographic Search Results

The initial database search using the keyword criteria returned many results that need to be filtered to only include works that are relevant to the four research questions outlined in Section 3.1. The filtering process is described using the PRISMA flowchart depicted in Fig. 1. Common examples of papers that were filtered from these results include the use of unrelated methods (e.g. natural language processing, blockchain, internet of things) or unrelated EHR modalities (e.g. images, diagnostic measurements, or unstructured text).

Through the initial database search, we identified 1939 empirical studies as candidates for review. Among those, 917 studies were qualitatively selected as relevant studies, based on the study quality assessment and exclusion criteria. The same steps were applied to the 42 studies identified through snowballing and we identified 1 additional relevant paper for our study. One additional study was included based on expert recommendations. To avoid the inclusion of duplicate studies which would inevitably bias the result of the synthesis, we thoroughly checked if very similar studies

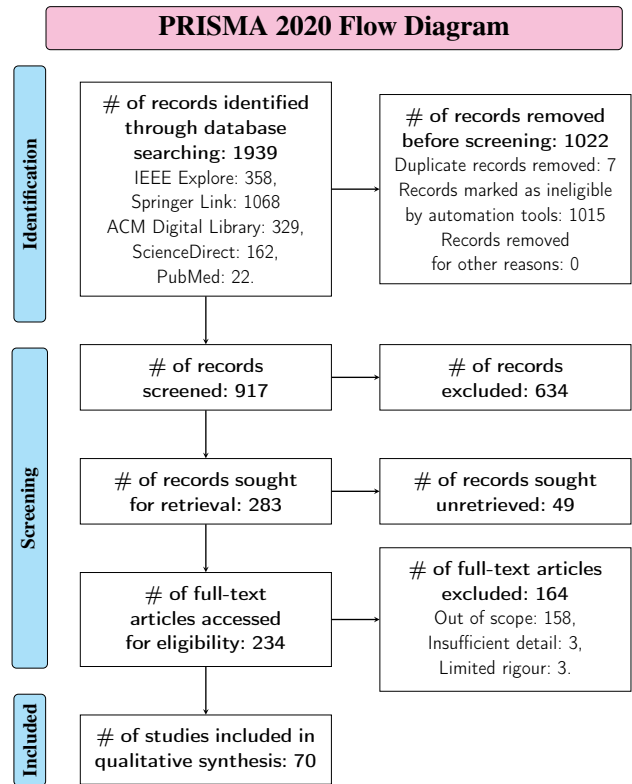


Figure 1: Illustration of the PRISMA filtering process [22] used on bibliographic search results.

were published in more than one paper. After assessing 234 full-texts, we eventually selected a total of 70 studies to be included in the synthesis of evidence.

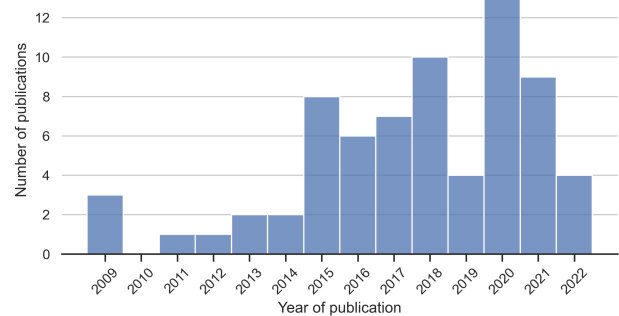


Figure 2: Histogram depicting the frequency of papers published by year for the 70 studies used in this systematic review. This figure shows that research in this area is increasing over time, beginning around 2011. Note that year 2022 does not contain the full-year of publications.

3.4. Exclusion Criteria and Quality Assessment

All the solutions identified from the search phase were reviewed for relevancy. Below are the exclusion criteria we adapted from Khan et al. [23] to screen studies of interest.

- Books and news articles, vision papers, or front-matter

- Papers without any data analysis (e.g system architecture for data management, or data pre-processing methodologies)
- Process mining papers (see Section 2.3.2 for justification)
- Papers that did not exclusively analyse administrative data, or whose data did not fit the characteristics of AHRs defined in Section 2.2.1
- Papers not written in English
- Papers whose full-text was not available for public access or through digital library services
- Papers published more than 10 years ago²

To ensure the quality of the search process, the initial extraction of 42 studies and preliminary analysis were done by manual web search with a combination of different keywords, snowballing and reverse snowballing of identified papers. Two authors went through the findings and thoroughly reviewed them to define the scope of our study and search keywords for the systematic review.

Once a formal literature search was conducted, one author screened titles and abstracts of 917 records, and assessed the full-text of 234 records. In each case, the author classified the relevancy of each study as ‘yes’, ‘no’, and ‘maybe’. At least two authors reviewed the ‘maybe’ studies and collectively decided to include or exclude the study from the evidence. Once the evidence was finalised, one author extracted results from the identified literature and a second author reviewed and verified the results to ensure accuracy.

4. Results

This section first presents an overview (Section 4.1) of the 70 studies identified through the bibliographic search process, before proceeding to use these studies to answer the four research questions (Sections 4.2-4.5).

4.1. Overview of AHR-based Studies

Fig. 2 depicts the year of publication for each of the AHR-based studies, and shows that AHR-based research is increasing over time. Among these studies, the following were the most cited by other AHR-based studies:

- With nine citations: [24] explores methods for representation learning of medical and diagnosis codes
- With eight citations: [25] uses a novel topic-modelling approach to mine latent treatment patterns in AHRs.
- With five citations: [26] identifies temporal disease trajectories using a population-scale AHR database.

²Some studies published more than 10 years ago may still be included if they were selected as a part of the informal keyword search.

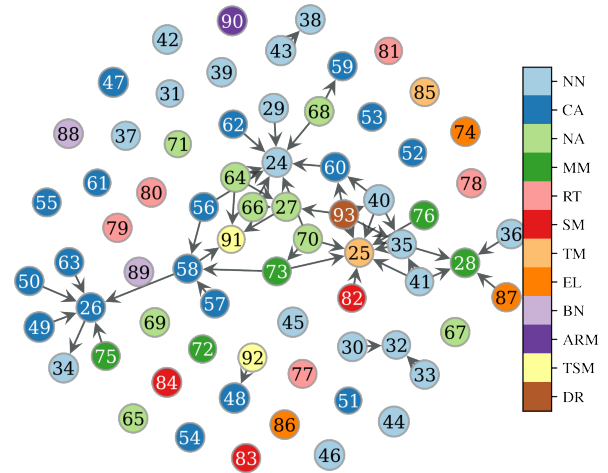


Figure 3: This citation network of the 70 AHR-based studies reviewed in this paper. Each node in the network is one study (numbers refer to references as per the bibliography), with directed edges indicating a *from*→*to* citation, and with node colour indicating the ML technique used in the study (as per Table 1).

- With four citations each: [27] uses a graph-based framework for temporal phenotyping, [28] defines clinical pathways observed in chemotherapy, and [29] explores the use of healthcare representation learning using graph-based attention models.

Fig. 3 depicts the citation network of the AHR-based studies, with each study coloured by the ML technique used (Section 4.2 explores these ML techniques in more detail). The mean number of AHR-based citations per study is 0.8, with a majority (53%) of studies containing no citations to other AHR-based literature. Furthermore, while many neural network-based studies cite each other, there is no clear technique-based citation patterns for other methods used. These findings indicate that AHR-based research is highly disconnected.

4.2. RQ1: What are the machine learning techniques utilised in analysing health service patterns?

This section explores which machine learning techniques are used to analyse AHRs. We provide a detailed analysis for each of the most common techniques, and a shorter analysis for some of the infrequently used techniques. The detailed analysis includes a flow-chart of the steps followed by the studies in applying each technique³ describing how each of the studies implement these analytical steps, and a discussion that outlines any other interesting findings from the studies.

In total, we observed twelve types of ML techniques used in the identified studies as summarised in Table 1.

³Evidence only report best performing models in identified studies and do not report models developed for comparison such as baseline methods.

Table 1
AHR analysis techniques

Technique	Studies
Neural Networks (NN)	[24, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46]
Cluster Analysis (CA)	[26, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63]
Network Analysis (NA)	[27, 46, 49, 50, 51, 59, 64, 65, 66, 67, 68, 69, 70, 71]
Markov Model (MM)	[28, 48, 72, 73, 74, 75, 76]
Regression-based Techniques (RT)	[66, 77, 78, 79, 80, 81]
Sequence Mining (SM)	[68, 69, 82, 83, 84]
Topic Modeling (TM)	[25, 39, 56, 85]
Ensemble Learning (EL)	[74, 86, 87]
Bayesian Networks (BN)	[75, 88, 89]
Association Rule Mining (ARM)	[84, 90]
Temporal Signature Mining (TSM)	[91, 92]
Dimensionality Reduction (DR)	[93]

Sections 4.2.1 to 4.2.5 include a detailed analysis of the five most common techniques, and Section 4.2.6 provides shorter analysis of each of the remaining techniques.

4.2.1. Neural Networks

Neural networks are the most widely adopted machine learning technique with nineteen studies. The analysis in these studies follows the four analysis steps as in Fig. 4. Specific techniques and details associated with each step are listed in Table 2, with associated studies for each technique.

When analysing the *Class of Techniques* row of Table 2, it is clear that the main contribution proposed in neural network-based studies is a representation learning strategy to capture and summarise qualities of AHR data and to reduce dimensions of AHRs. Five studies propose end-to-end frameworks that encapsulate both representation learning and predictive learning.

In the feature preparation step, different features or combinations of features from EHR records are used with diagnosis and treatment codes used more prominently than others. Two studies go beyond common feature representation and incorporate domain knowledge to increase representation quality. In Li et al. [35] and Ochoa and Mustafa [46], instead of using medical records directly, vectorised patient representations resulting from auto-encoders were used as the feature for sequence embedding and graph embedding tasks respectively.

Under the learning step, different neural network-based models are used with sequence-based based frameworks (e.g., LSTM, RNN, Transformer), auto-encoders, and sequence embedding being the most commonly applied approaches. Many studies propose novel techniques or frameworks for feature representation as well.

As most studies focus on representation learning, the main usage of the neural network model in these studies is for some downstream analytics tasks such as prediction or clustering. Two studies each use the dimensionality reduction provided by representation learning to conduct visual analytics and to generate clinically meaningful interpretations of

large datasets. All ten studies evaluate their proposed solution by comparing performance, usually by conducting one or more prediction tasks and comparing performance with baseline methods for representation learning and predictive analytics.

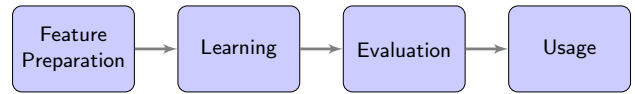


Figure 4: Steps followed in neural network-based studies

4.2.2. Cluster Analysis

Cluster analysis is another widely used technique with eighteen studies utilising multiple different similarity measures and clustering algorithms. Fig. 5 illustrates the generic process steps followed in these studies. We observed that the first three steps (data representation, calculating similarity, and clustering) are followed in all studies, but some studies only followed one of the interpretation and evaluation steps.

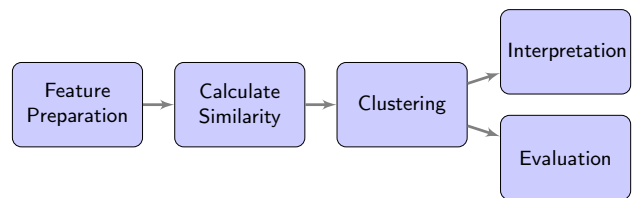


Figure 5: Steps followed in cluster analysis-based studies

Table 3 catalogues the different ways cluster analysis is used in the identified studies.

The most widely used data representation technique is to represent patient clinical pathways as strings with different events encoded as letters. In addition to that, a variety of vector representations are used in different clustering algorithms for different purposes, including the novel model by Chen et al. [57] that represents the drug use duration statistics of a patient as a vector. Finally, two studies [51, 59] utilised networks to represent the data.

Table 2
Neural Network techniques

Step	Technique	Studies
Class of Techniques	Representation Learning	[24, 29, 30, 31, 32, 33, 34, 35, 36, 37] [38, 40, 42, 44]
	Predictive Learning	[32, 34, 39, 43, 45, 46]
Features Used	Diagnosis Codes	[24, 29, 30, 33, 34, 38, 39, 40, 42, 43, 44, 45, 46]
	Treatment Information (procedure and medication codes)	[24, 30, 33, 34, 35, 38, 39, 40, 42, 43]
	Encounter Records	[30, 31, 32]
	Domain Knowledge (grouping)	[29, 30]
	Demographic Information	[34, 35, 46]
	Vectorised Patient Representation	[35]
Learning Technique	LSTM based framework	[32, 34, 37]
	Autoencoder	[31, 37]
	Sequence Embedding	[31, 35, 38]
	RNN based framework	[35]
	Reinforcement learning	[43]
	Graph Neural Network	[46]
	Multilevel Embedding (Novel)	[30]
	GGraph-based Attention Model (novel)	[29]
	Med2Vec (novel)	[24]
	Patient2Vec (Novel)	[33]
	RoMCP (novel)	[36]
	Neural Topic Model (Novel)	[39]
	ProAID (Novel)	[40]
	Clinical Language Model-based Representation (Novel)	[42]
	Categorical Sequence Encoder (CaSE) (Novel)	[44]
Multi-task Transformer (TransMT) (Novel)	[45]	
Usage	Feature for downstream analytics	[24, 29, 30, 32, 33, 34, 35, 36, 37, 38] [39, 40, 42, 44]
	Visual Analytics	[31, 33]
	Clinically meaningful interpretations	[24, 33, 43, 45, 46]
Evaluation	Compare performance with baseline methods	[24, 29, 30, 31, 32, 33, 34, 35, 36, 37] [38, 39, 40, 42, 43, 44, 45, 46]
	Visualise and inspect clusters	[46]

To calculate similarity between events in AHRs for cluster analysis, most studies use common distance measures such as Euclidean distance, Jaccard score and Levenshtein distance, but the choice of metric depended on the data representation and the study objective. We also identified five studies that have proposed novel similarity measures (or variations on the most common measures) that can capture unique properties of health service patterns.

Multiple different clustering algorithms were utilised under the clustering step. Ward's method for hierarchical clustering was the most commonly used algorithm. Hierarchical-based clustering algorithms were also common, particularly because of the flexibility it offers to adjust the number and size of resulting clusters.

Interpretation of clustering results was commonly achieved using descriptive analytics that incorporates visualisation techniques including heat-maps, dendrograms, and charts of statistical properties. Additionally, four studies extracted cluster cores or representative samples to interpret the characteristics of the clusters. For example, Chen et al. [52] extracted cluster cores using KNN and Zhang et al. [48] constructed a transition matrix of Markov chains to identify

and connect patient states that occur with high probability and frequency, within a cluster.

Evaluation of clustering results was commonly achieved by comparing cluster quality (e.g. inter- and intra-cluster distance, mutual information, or Rand index) against clusters generated via benchmark methods, and by verifying the meaningfulness of the clusters against external information or domain knowledge (such as patient demographics and clinical best practices). Another evaluation technique was to use cluster membership as a feature in the classification algorithm, using it to optimise the classification accuracy of data and to indicate that meaningful clusters were formed during the analysis.

4.2.3. Network Analysis

A total of fourteen studies use network analysis to represent and analyse health service patterns. Network analysis based studies were observed to follow five steps as shown in Fig. 6. The specific techniques and studies associated with each step are listed in Table 4.

In network analysis studies, nodes most often represent diseases (six studies) or clinical events (five studies), while

Table 3
Cluster analysis techniques

Step	Technique	Studies
Data Representation	Clinical pathway as a string	[48, 53, 55, 56, 58]
	Vector of diagnosis codes by occurrence frequency	[47, 49]
	Treatment record set sequence	[52, 93]
	Vector of patient state	[54, 60]
	Disease co-occurrence vector	[50]
	Diagnosed disease trajectory	[26]
	Cancer metastasis dynamic network	[51]
	Patient flow dynamic network	[59]
Drug use duration statistic vector (novel)	[57]	
Similarity Measure	Euclidean distance	[47, 57, 59]
	Jaccard score	[26, 50]
	Cosine similarity	[49, 60]
	Levenshtin distance	[53]
	Damerau–Levenshtein distance	[56]
	Metastasis conditional incidence (hazard) functions	[51]
	Longest common subsequence	[48]
	Maximal Repeat Alphabet Feature Set (novel)	[55]
	Fusion of 3 novel similarity measures- content sequence and duration view (novel)	[93]
	Modified Needleman–Wunsch algorithm (novel)	[58]
	Markov theory based measure (novel)	[52]
Ordinal edit distance (novel)	[54]	
Clustering Algorithm	Hierarchical agglomerative clustering with Ward's method	[33, 47, 50, 51, 54, 55, 56]
	Hierarchical clustering based on complete linkage	[53]
	Hierarchical clustering based on average linkage	[49]
	Affinity propagation	[52, 57]
	Markov Cluster Algorithm	[26, 60]
	Spectral Clustering	[93]
	K-means	[93]
	K-medoids clustering	[58]
PCA based clustering	[59]	
Interpretation	Descriptive analytics and visualisations	[47, 49, 51, 53, 54, 57, 59, 60, 93]
	Extract cluster core/representative samples	[48, 52, 53, 56]
	Compare with clinical guidelines	[93]
Evaluation	Cluster quality	[54, 55, 57, 58, 93]
	Verifying interpretability of clusters	[26, 48, 52, 57, 60]
	Using cluster membership as a feature for classification	[50, 93]

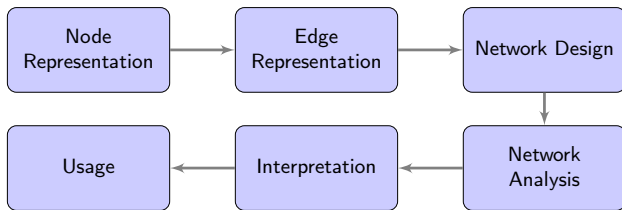


Figure 6: Steps followed in network analysis-based studies

edges most often represent co-occurrence of records (six studies) or a temporal relationship (five studies). Certain studies have reduced the network dimensions by means such as grouping the diseases in co-morbidity networks using an ontological hierarchy of ICD-9-CM codes. Co-occurrence measures used to define relationships in networks varied from simple measures such as co-occurrence frequency to

TF-IDF and statistical equations such as pairwise association analysis using a chi-square test and correlation. We observed that all network designs are weighted, with five directed networks. Three studies each have utilised temporal and dynamic networks designs as well. Two studies [51, 59] were observed to utilise unique network representations to look into cancer metastasis spread and patient flow in hospital wards.

Various techniques have been utilised to analyse and interpret the resulting network as listed in Table 4, with a majority of studies using visualisations and statistical interpretations of network properties. Some studies have utilised network properties such as hubs, dense cores, communities and degree variability to understand and interpret network behaviour.

Network representations were used widely for descriptive analytics and to generate measures for downstream analysis

Table 4
Network analysis techniques

Step	Technique	Studies
Node Representation	Disease (Co-morbidity network)	[49, 50, 64, 65, 66, 67]
	Patient	[46, 49, 64]
	Clinical event	[27, 68, 69, 70, 71]
	Sites of metastasis	[51]
	Hospital wards	[59]
Edge Representation	Measure of co-occurrence in records	[46, 49, 50, 64, 65, 68, 71]
	Temporal relationship	[27, 59, 67, 69, 70]
	Risk of developing subsequent disease	[51, 64, 66]
	Causal information fraction	[67]
Network Design	Weighted	[27, 46, 49, 50, 51, 59, 64, 65, 66, 67, 68]
	Directed	[27, 51, 59, 65, 66, 67, 68, 69]
	Temporal	[51, 59, 65, 70]
	Dynamic	[51, 59, 64]
Interpretation	Visualisation	[27, 49, 51, 65, 66, 67, 69, 70]
	Statistical properties of network	[49, 51, 64, 66, 67, 71]
	Hubs	[64, 66]
	Dense core	[49, 64]
	Community detection	[64, 67]
	Long-range edges	[64]
	Degree stability over time	[59]
	Edge weight variability over time	[59]
	Extract temporal phenotypes (representative network)	[27]
Random walks	[68]	
Usage	Descriptive analytics	[27, 49, 51, 59, 64, 65, 66, 69, 71]
	Generate measures for downstream analysis	[27, 46, 49, 50, 51, 64, 67]
	Pattern detection	[27, 66, 67, 68, 69]
	Observe emergent behaviour over time	[64]

tasks such as clustering or risk prediction. Descriptive analysis, pattern detection and observing emergent behaviour over dynamic networks were also objectives of some studies.

A noteworthy observation is that network analysis studies often do not include an evaluation step. We attribute this to the inherent explainability of network models and simple process of network representation and design, which makes the network analysis results intuitive, traceable and any interpretation and usage easily explainable.

4.2.4. Markov Models

We identified seven studies that follow a Markov model-based analysis for health service pattern discovery in AHR. Markov model-based studies follow five steps as illustrated in Fig. 7, and Table 5 lists the techniques applied in each step.

The most common way that state is represented in the Markov modelling-based studies is as either patient health state or a sequence of events and dynamic features. The Markov models applied in the studies are often constructed using more complex variations of the Markov model, such as a hidden Markov model or other novel extensions. This is to account for some hidden complexity that isn't captured in AHRs [74] or to incorporate some domain knowledge into the model [73].

Evaluation is based on domain knowledge and expert consultation, comparing performance with baseline methods and by numerical simulations. Probability-based descriptive analytics and visualisation are used to interpret Markov model results, but are used for different objectives such as understanding treatment pathways, for preprocessing and interpretation, and for generating features for downstream analytics.

It is important to note that healthcare is inherently a non-Markovian process, since the events recorded in AHR datasets contain long-term dependencies and do not explicitly depend on the previous state. For example, a routine admission with irrelevant medical information would destroy the illness memory, especially for chronic conditions [32]. This should be considered carefully when applying Markov models for health service pattern discovery.



Figure 7: Steps followed in Markov-based studies

Table 5
Markov Modelling techniques

Step	Technique	Studies
State Representation	Patient health state	[28, 72, 76]
	Sequences of events and dynamic features	[48, 73, 74]
Modelling Technique	Markov chain	[48]
	Finite horizon markov decision process	[72]
	Markov model	[28]
	Hidden Markov Model (HMM)	[74, 76]
	Novel extension of HMM	[73]
Evaluation	Domain knowledge and expert consultation	[28, 48, 76]
	Compare performance with baseline methods	[73, 74]
	Numerical simulations	[72]
Interpretation	Probability based descriptive analysis	[28, 48, 72, 73]
	Visualisation	[48, 72, 73]
Usage	Understand and improve treatment pathways	[28, 72, 73]
	Preprocessing and interpretation	[48]
	Feature for downstream analytics	[74]
	Correlated with medical-oriented outcomes	[76]

4.2.5. Regression-based Techniques

We identified six studies that use regression-based techniques to analyse health service patterns, following the five steps illustrated in Fig. 8. Table 6 lists the techniques used in different studies for the four steps followed by variable selection.

All of the regression-based studies use similar variable characteristics to select which variables would be used in the analysis, such as whether they are dependent/independent, and which indexes or metrics are used to measure them.

The most common regression model used was multivariate logistic regression, but other variations were also observed. Similarly, all of the regression-based studies evaluate their findings using statistical significance and sensitivity analysis, as well as using other evaluation techniques and validating results with domain knowledge where applicable [66]. Te Marvelde et al. [78] used an additional validation step to ensure survival difference is independent of the age-at-diagnosis by using the multivariable Cox proportional hazard method.

All studies use descriptive statistical analytics and/or visualisations to interpret their results. Finally, each of the studies use the regression analysis to better understand clinical practice and associations between AHR variables. Glicksberg et al. [66] quantifies the risk associated with diseases by using their regression analysis as a preprocessing step for edge weights in a disease network.

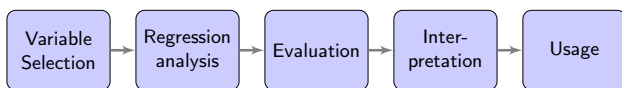


Figure 8: Steps followed in regression-based studies

4.2.6. Other Techniques

This section summarises the studies that utilise less common analytics methods. We briefly describe how these

techniques are applied according to the features used, modelling techniques, and, where it was observed, methods for evaluation, interpretation, and usage. Our synthesis of the approaches to using these techniques may be limited in some cases due to the modest number of studies that utilise these methods.

Sequence Mining We identified five studies used sequence-based mining approaches. In these studies, authors cited the importance of capturing relationships between the order of events within AHRs. For example, Kaur et al. [82] proposed a time range based sequence mining approach on treatment data to generate informative features for downstream analysis. The methods of these studies are further detailed in Table 7.

Topic Modelling We identified four studies that use natural language processing-inspired topic modelling techniques (Latent Dirichlet Allocation (LDA) and Additive Regularisation of Topic Models (ARTM)), and one novel technique, Treatment Pattern Model (TPM). These studies identify health service patterns and use them for understanding clinical pathways. A common theme of these studies is the modelling of latent variables that are indirectly associated to the recorded variables in AHRs. The methods of these studies are further detailed in Table 8.

Ensemble Learning We identified three studies that use some form of ensemble-based learning techniques. These studies each pursue prediction-based tasks, often using multiple AHR attributes. The methods of these studies are further detailed in Table 9.

Table 6
Regression techniques

Step	Technique	Studies
Analysis Technique	Multivariate logistic regression	[66, 77, 79, 81]
	Ordered logistic regression	[78]
	Proportional hazards regression	[77]
Evaluation	Statistical significance and sensitivity analysis	[66, 77, 78, 79, 81]
	Confidence level	[77, 78, 79]
	Multivariate Cox proportional hazard method	[78]
	Against known data and clinical expectations	[66]
Interpretation	Descriptive statistical analytics with tables	[66, 77, 78, 79, 81]
	Visualisation	[66, 78, 81]
Usage	Discuss clinical implications	[66, 77, 79, 81]
	Identify new associations	[77, 78, 79]
	Preprocessing step	[66]

Table 7
Sequence Mining

Step	Technique	Studies
Features Used	Treatment data	[82, 83, 84]
Modelling Technique	Time range based sequence mining	[82]
	Transitive Sequential Pattern Mining (novel)	[83]
	Sequential Relational n-Disease Pattern (SERENDIP) (novel)	[84]
Evaluation	Compare prediction accuracy with baseline methods	[82, 83]
	Prediction accuracy	[84]
Interpretation	Identify significance of various attributes	[82, 83, 84]
Usage	Features for downstream analytics	[82, 83]

Table 8
Topic Modelling

Step	Technique	Studies
Features Used	Patient features and treatment behaviors	[25, 85]
	Pathway as a text string	[56]
Modelling Technique	Latent Dirichlet Allocation (LDA)	[25]
	Additive Regularisation of Topic Models (ARTM)	[56]
	Treatment Pattern Model (TPM) (novel)	[85]
Evaluation	Compare performance with baseline methods	[25]
	Perplexity score and sparsity for optimal topic numbers	[56]
Interpretation	Visualisation	[25]
	Probability based descriptive analysis	[56]
	Discover underlying treatment patterns	[85]
Usage	Understand clinical pathways	[25, 56, 85]

Table 9
Ensemble Learning

Step	Technique	Studies
Features Used	Clinical events	[86, 87]
	Administrative events	[41, 86, 87]
	Patient demographics	[41, 87]
Modelling Technique	Random Forest	[86]
	Stochastic Gradient Boosting (GBM)	[41, 87]
	Multiple logistic regression	[41, 87]
Evaluation	Prediction accuracy	[41, 86, 87]
	Compare performance with baseline methods	[41]
Interpretation	Measure phenotype severity	[86]
	Predict psychiatric hospital care	[41]
	Predict hospital readmission	[87]
Usage	Outcome assessment	[86]
	Decision support	[41, 87]

Table 10
Bayesian Networks

Step	Technique	Studies
Features Used	Patient health state	[75]
	Demographic, diagnosed-based and prior-utilization variable	[88]
	EHR events corresponding to patient risk factors	[89]
Modelling Technique	Bayesian Networks	[75, 88]
	Multiplicative-Forest Point Processes (MFPPs) (novel)	[89]
Evaluation	Prediction accuracy	[75, 88, 89]
	Compare performance with baseline methods	[88]
Interpretation	Probability density function based	[88]
	Likelihood of disease	[75, 89]
	Disease progression trajectory	[75]
Usage	Identify disease risk factors	[88]
	Predict survivability	[88, 89]
	Improve understanding of diagnosis severity	[75]

Table 11
Association Rule Mining

Step	Technique	Studies
Features Used	Clinical Events	[84]
	Diagnosis codes	[92]
Modelling Technique	Sequential, Relational, n-Disease Pattern (novel)	[84]
	A priori algorithm	[92]
Interpretation	Frequent item sets	[84]
	Pattern discovery	[92]
Usage	Multi-morbidity disease prediction	[84]
	Discover toxicity patterns	[92]

Bayesian Networks Bayesian Networks were used in three studies to identify risk factors, and to predict survivability and progression of diseases. The probability-based interpretation of the model was cited as a particular advantage of Bayesian Networks. The methods of these studies are further detailed in Table 10.

Association Rule Mining Association rule mining was used on only two identified studies, typically for high-level AHR analysis – such as pattern discovery and identifying frequent item sets – and are further detailed in Table 11.

Temporal Signature Mining We identified two studies proposed temporal signature mining, which is distinct from sequence mining, since the time between events is considered in addition to the sequence-based nature of the data. Wang et al. [91] represent and analyse patient encounters through a novel matrix representation, and their approach is further detailed in Table 12.

Dimensionality Reduction Only one study [93] used Laplacian Eigenmaps (LEs) to embed the adjacency graph of a large network into a low-dimensional space to

Table 12
Temporal Signature Mining

Step	Technique	Studies
Features Used	Patient encounters	[91]
	Diagnosis codes	[92]
Modelling Technique	Novel Matrix Representation	[91]
	Extension of convolutional non-negative matrix factorisation	[91]
	A priori algorithm	[92]
Evaluation	Clinical interpretations of identified signatures	[91]
	Performance in detecting known patterns in synthetic data	[91]
Interpretation	Visualisation	[91]
	Clinical significance	[91]
	Pattern discovery	[92]
Usage	Understand latent event patterns of clinical significance	[91]
	Discover toxicity patterns	[92]

Table 13
Dimensionality Reduction

Step	Technique	Studies
Features Used	Multiple similarity networks from treatment records	[93]
Modelling Technique	Laplacian Eigenmaps	[93]
Usage	For downstream analysis	[93]

conduct features for downstream analysis. See Table 13. Note that other approaches – such as neural networks for representation learning (Section 4.2.1) or vector-based representations used for clustering (Section 4.2.2) – may also effectively reduce the dimensionality of data. However, these are not included here since dimensionality reduction is not the primary motivation for such methods, nor was it the stated objective of the authors in such studies.

4.3. RQ2: What are the applications of health service patterns in the identified studies?

We mapped the identified literature into eight key application areas that use machine learning techniques on EHR data. These applications are adapted from other applications recognised by existing EHR-based mining analytics surveys [3, 4], however we focus the scope of these applications specifically to AHR-based HSPM.

Table 14 shows which applications are pursued in each of the reviewed studies. We found that the most common applications were analysis of healthcare patterns and of medical trajectories. We believe this is for two reasons: 1. These applications are valuable to health informatics research, and 2. The characteristics of AHRs (e.g. their temporal nature) make these applications practical.

Other common applications included: comorbidity analysis, healthcare guidelines, cohort identification, and risk prediction. These studies often use very distinct or unambiguous features of AHRs, such as diagnosis codes to predict risk [70, 88] or identify comorbidities [47, 50]. Finally, only three studies pursue outlier detection, while only one study pursues intervention analysis. This may be due to a lack of distinct or unambiguous features in AHRs that correlate to high-level healthcare concepts such as interventions and anomalies.

The remainder of this section describes each of these applications in some more detail.

Healthcare patterns

The most common application area was to understand the patterns and pathways in medical care that patients receive. This may include patterns of drug or procedure codes that are commonly prescribed (such as learned representations of medical-codes [24]), or care pathways across a healthcare system [48]. In total, 27 studies attempted some form of healthcare pattern analysis.

Medical trajectory

Another common application, observed in 23 studies, was prediction of future medical events in the treatment

sequence of a patient and prediction of other future phenomena such as outcome or readmission. For example, Jensen et al. [26] use large, population-wide data registry to extract temporal disease trajectories.

Comorbidity analysis

This application is concerned with identifying patients that express two or more diseases or medical conditions at the same time. In one study, Sideris et al. [50] develop a feature-extraction framework for identifying comorbidities in AHRs. We identified ten studies in total that applied AHR-based analytics for comorbidity analysis.

Healthcare guidelines

Nine studies sought to use AHRs to inform the development of evidence-based guidelines in healthcare. For example, Roder et al. [77] observes correlations between country-of-birth and female breast cancer in New South Wales, Australia to draw implications for health-service delivery.

Cohort identification

Eight studies identify similar groups of patients within an AHR database according to an attribute or set of attributes. For example, Roque et al. [49] identify patient cohorts within a population using a hierarchical stratification.

Risk prediction

Seven studies quantify patient severity or risk based on attributes contained within AHRs. For example, Wang et al. [88] determines a survivability prognosis for lung cancer patients by analysing relevant risk factors.

Outlier detection

This application relates to identifying patients who receive care that is not in line with that of similar patients. For example, Hompes et al. [60] identifies variations in healthcare processes in event data from AHRs.

Intervention analysis

The final observed application area relates to evaluating the efficacy of healthcare interventions. In the one study that pursued this application, Te Marvelde et al. [78] used existing care pathways to observe survival outcomes among colon cancer patients.

4.4. RQ3: What are the analytics activities associated with health service pattern mining?

In sections 4.2 and 4.3, we highlight the twelve machine learning techniques and eight HSPM applications identified

Table 14
Applications of Health Service Patterns

Application	Studies
Healthcare patterns	[24, 25, 27, 28, 29, 31, 33, 36, 38, 40, 42, 44, 46, 48, 52, 54, 55, 56, 61, 64, 68, 69, 73, 79, 83, 91, 93]
Medical trajectory	[26, 27, 30, 32, 35, 37, 39, 40, 41, 43, 45, 51, 62, 64, 70, 72, 74, 75, 80, 82, 85, 87, 89]
Comorbidity analysis	[47, 49, 50, 62, 64, 66, 67, 76, 84, 84]
Healthcare guidelines	[34, 52, 57, 58, 59, 61, 63, 77, 90]
Cohort identification	[24, 49, 50, 53, 65, 71, 86, 92]
Risk prediction	[70, 71, 75, 81, 86, 88, 92]
Outlier detection	[60, 80]
Intervention analysis	[78]

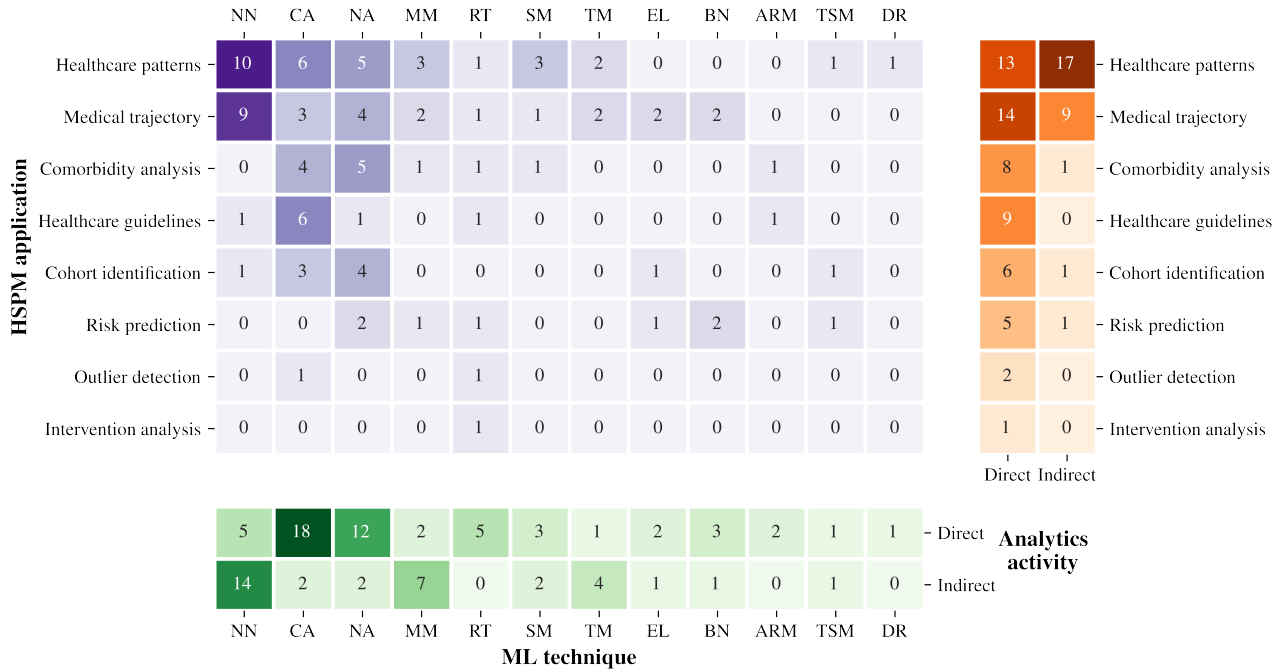


Figure 9: A 2D heatmap (purple) depicting how many studies (cell-value) utilise a given machine learning technique for each HSPM application. Also pictured are which studies undertook direct/indirect analytics activities, by both techniques (green) and applications (orange).

in literature. This section looks at how these techniques are used in pursuit of these applications. First, we look at how often each technique is used for each application. Then, we look at the analytics activity undertaken in each study to gain further insight into the research.

The 2D heatmap in Figure 9 depicts the distribution of ML techniques and HSPM applications. It shows three major technique-related trends.

- Clustering and network analysis are applied often and in many and diverse application areas, including healthcare patterns, medical trajectories, comorbidity analysis, healthcare guidelines, and cohort identification. A similar relationship was observed for network analysis, except it was used less often for healthcare guidelines and more often for risk prediction.
- Neural networks were applied almost exclusively to healthcare patterns and medical trajectory analysis.

Also, Markov models, topic modelling, ensemble learning, Bayesian networks, sequence mining, and dimensionality reduction were also used almost exclusively for these same applications, though they were not as commonly used as neural networks.

- There was no clear application-related trends for regression analysis, association rule mining, or temporal signature mining.

Furthermore, we also observed that the majority of studies (64%) often use the associated machine learning techniques to directly pursue their desired application. A common example of a direct analytics activity is analysis of clusters or networks of drug or diagnosis codes. However, a significant number of studies (36%) model healthcare attributes that are not recorded in AHRs indirectly using what is available in the data. For example, a common analytics

activity is to learn more useful representations of medical codes that reveal latent medical- or healthcare-knowledge [24, 29, 30]. To distinguish between these two types of analytics activities, we refer to them as direct-analytics and indirect-analytics respectively, and define them as follows:

Direct-analytics Modelling or prediction of attributes that have clear definitions and observations.

Indirect-analytics Modelling or prediction of attributes that are either conceptually abstract or are not directly observed in the data.

We distinguished each study identified in this review using these definitions, and present the results in Figure 9 adjacent to the primary heatmap depicted. When distinguishing the techniques in this way, we find that the following techniques almost exclusively pursue direct-analytics tasks: clustering, network analysis, regression analysis, and association rule mining. In contrast, neural networks, Markov modelling, and topic modelling are used almost exclusively for indirect-analytics tasks. Furthermore, when distinguishing applications in this way, we find that indirect-analytics tasks are common for analysing healthcare patterns and medical trajectories, while direct-analytics tasks were common among every application.

These findings suggest the choice of machine learning technique used should be informed by the analytics task. Furthermore, indirect-analytics is more applicable to the modelling of the latent and often abstract attributes not recorded in AHRs, and is common in modelling healthcare patterns and medical trajectories.

4.5. RQ4: What are the limitations of existing machine learning techniques for health service pattern discovery?

This section discusses several important limitations regarding AHR and machine learning techniques for health service pattern discovery.

4.5.1. Limitations Associated with Data

This subsection outlines seven relevant limitations that pertain to AHRs.

Data Quality

A common limitation we identify in literature is that of data quality, which can affect the performance, validity, and applicability of particular modelling approaches. Typical data quality issues include: missing values, incorrect values, inconsistent formatting, inconsistent reporting standards, or noise introduced into data from human or instrument errors.

Some machine learning techniques may be robust to particular data quality issues, however their applicability in such cases depends on the type and severity of the issues. For example, clustering was often applied to noisy datasets, but not for datasets with missing values. This is because many similarity measures used in clustering algorithms are tolerant of small variations in the data, but cannot compute

similarity where data is missing [54]. Other means of addressing data quality issues specific to AHRs include using ensemble-learning approaches and imputation algorithms for handling missing values [94, 95] and class imbalance [96].

Data Accessibility

AHRs contain information pertaining to a patient's health state and health service access, as well as personally identifiable attributes such as age, sex, or geographical locality. Access to AHRs is often restricted to limit the risk of the sensitive information being leaked or breached [97]. In a research context, these restrictions have implications for what parts of AHR-based research can be shared, what findings can be revealed, and limit the ease with which research findings can be reproduced [98].

Indirect Nature of Data

AHRs do not specifically seek to make an accurate recording of a patient's health state. The information recorded within AHRs is still valuable for health informatics since it indirectly relates to a patient's health state. For example, the primary purpose of billing codes recorded by hospitals or insurance providers is accounting, but these codes may be used to indirectly infer a patient's health state since they indicate which services the patient received.

The key limitation of indirect data is noise and uncertainty, and it is critical to understand and quantify uncertainty when using indirect data. Using uncertain quantities as an evaluation criterion may lead to misleading or inaccurate research conclusions. For example, many of the studies reviewed in this paper explore the relationship between patient readmission (or lack thereof) and existing treatment patterns. In these studies, patient readmission is used as a proxy for the success or failure of a treatment. However, this is an indirect attribute for treatment success, since there are many reasons why patients may/may not be readmitted.

Researchers should practise caution when using indirect relationships within AHRs to make causal assertions or predictions and, where possible, compute confidence intervals and estimate likelihood of any conclusions drawn from indirect means.

Fragmentation

The most accessible types of AHRs are those sourced from a single institution, such as hospital administration data, insurance claims data. Analysis of single-source datasets exposes the indirect- and incomplete-data quality issues discussed above. To move beyond single-source AHRs, patients need to be linked across the various disparate single-source AHRs. We refer to this issue as AHR fragmentation. While many methods exist for linking entities across disparate databases, the aforementioned issue of data accessibility makes this difficult. The process of linking data further accentuates the data accessibility limitation, but linked databases may mitigate some data quality or completeness limitations.

Single-source AHRs are often not sufficient for population-level health studies, since the relevant patient records may be distributed amongst multiple healthcare providers or across jurisdictional borders. The data fragmentation limitation significantly inhibits health informatics research at a population-level.

Irregular Temporal Data

Recordings in AHRs are typically not done in a structured manner. For example, two patients with same diagnosis and treatment patterns may have distinct AHRs, since they may not follow the same visit frequency, or their institution or jurisdiction may observe different administrative recording practices. This is a significant limitation in AHR-based research, with some studies choosing to represent the temporal data simply as a sequence. Moreover, other studies – especially those that calculate patient similarity using Clustering – simply disregard temporal information entirely, aggregating information within discrete bins such as hospital visits instead. This indicates that the temporal information stored within AHRs is difficult to utilise, however not doing so risks losing valuable contextual insight from the data.

Capitalising on the temporal information requires researchers to come up with similarity measures or feature representations that rely on assumptions that are context specific, hard to validate and generalise. We identify there is a need to benchmark such temporal representations of AHR data and report their performance across different machine learning tasks.

Data Quantity

The quantity of data available in AHRs is another common limitation with regards to both the number of records (rows) and the number of attributes per record (columns). This problem is often a symptom of other limitations (such as fragmentation or data quality), but not always. For example, even population-level, comprehensively linked AHRs may not contain many records for patients with uncommon diseases. Insufficient data leads analytical methods to discover weak or inaccurate relationships or characteristics. Care should be taken to report on the confidence of predictions when AHRs are affected by this limitation to avoid weak confidences of results going unnoticed.

It may be possible in some cases to rectify this limitation to a certain degree with data augmentation, or with synthetic data. For example, Steinberg et al. [42] develop a model to learn representations of medical records such that it does not suffer significantly from limited data.

Incorporating domain knowledge

Domain-specific knowledge is often critical to making AHRs and AHR-based analytics interpretable, as well as for organising data, pre-processing data, and for feature engineering.

For example, we observe many studies use International Statistical Classification of Diseases and Related Health Problems (ICD) [99] codes. These codes are high-dimensional and not human-readable, though they have

human-readable definitions. Where feasible, these studies can benefit from ontologies – such as the Clinical Classification Software (CCS) [100] – to reduce the sparsity of such codes sets by incorporating semantic, hierarchical structure.

Aside from semantic structure, other code-grouping structures might also be clinically relevant. For example, specific treatments may involve grouping of semantically distinct codes, and combinations of specific procedure and diagnosis codes may indicate different degrees or severities of a particular disease (e.g., early-stage vs terminal cancer).

Another example of domain-specific knowledge is the use of clinical guidelines to inform the construction of features or the design of models. For example, Choi et al. [24] use clinical guidelines to inform the design of their model (e.g., grouping events by hospital admission), and Choi et al. [29] use clinical guidelines to inform the construction of features (e.g., hierarchical features via CCS).

Finally, we note that domain-specific knowledge is often difficult to incorporate into AHR-based analytics because the data is administrative and not in the healthcare domain. The semantic gap between the administrative and clinical domains is significant, and designing methods to bridge this gap (such as the CCS) requires significant domain expertise. Efforts to develop such domain-bridging methods can vastly increase the clinical utility of AHR-based machine learning analytics.

4.5.2. Limitations Associated with ML Techniques

This subsection outlines three relevant limitations that pertain to the ML techniques used.

Oversimplification of AHR data

Many of the AHR-based studies outlined explore simple relationships in AHRs, such as the similarity between drugs or diseases based on co-occurrence. The techniques used in these studies cannot capture complex health service patterns since the techniques used require simplification of temporal, highly-dimensional AHR data. When AHRs are simplified, much of their subtlety is lost, which precludes extraction of complex healthcare patterns and medical trajectories and pathways. Many studies acknowledge this issue, and seek to develop novel modelling techniques instead of simplifying the data. This trend is most evident with neural network-based techniques (Section 4.2.3), where the majority of studies use novel methods to capture such complex relationships, often through representation learning.

Lack of Population-scale Studies

Another limitation is that many studies limit the scope of their methods to patient- or hospital-level, and do not attempt to analyse population-level patterns. This trend can often be attributed to AHR-related limitations, however it is also evident that such population-level studies require significant abstraction, filtering, and assumptions informed by domain knowledge. This unique challenge is also not addressed by any established techniques, and instead requires the development of new and novel solutions.

Scalability

Memory requirements are a common limitation in ML analytics, particularly in network science and deep learning approaches. These methods often require a complex composition of data structures and algorithms, especially during model training. In many AHR-related studies, scalability issues are often mitigated by simplifying AHRs to the point that they lose much of the complex and temporal relationships inherent in the records.

5. Open Research Questions

This section discusses open research questions in the area of population-level health service pattern discovery we identified, based on the state-of-art literature and emerging research areas that utilise AHRs.

How can benchmark and synthetic data be used to accelerate AHR-related research?

Benchmark datasets are public datasets, and thus do not exhibit many of the limitations associated with AHR data, such as data accessibility, data quantity, and fragmentation. Benchmark datasets also make AHR-related research accessible to more researchers. Many of the studies explored in this review employ benchmark datasets for evaluating their methods. Despite widespread use of benchmark datasets, there are not many available to the public for research. The most commonly used benchmark datasets include the MIMIC-III [101] and -IV [102] (containing intensive care unit AHRs), and the BPIC-11 [103] (Gynaecology department AHRs). Benchmark datasets are critical to the development of novel modelling techniques, for evaluating and comparing these techniques with a wider audience, and for further progressing AHR-based research. Availability of many diverse benchmark datasets will help in these endeavours.

Data synthesis methods can also be a useful way to share datasets and develop new modelling techniques [104] as it can mitigate data privacy issues and data quality limitations. Many methods have been established specifically for synthesising patient-level data [105, 106]. In comparison to benchmark data, synthetic data is not a widely used technique for mitigating the various AHR-related limitations. Synthetic data should be used to help accelerate development of research in this area.

How can AHRs be accurately represented and visualised?

AHRs often contain indirect relationships to high-level healthcare concepts (as discussed in 2.2.1 and evidenced by the many indirect-analytics studies shown in Figure 9), which makes it challenging to explore the data to find these relationships. Many of the limitations associated with AHR-data – such as data quality, fragmentation, and irregular temporal data – also inhibit representation and visualisation efforts. This review highlighted many ways in which AHRs can be represented and visualised for domain-specific

tasks, however many of these representations and visualisations thereof are abstract. New methods could be developed to make interpretable representations and visualisations of AHRs.

Can dimensionality reduction be applied to temporal AHR data?

Some of the identified studies use dimensionality reduction as a data pre-processing or post-processing step, generating features for downstream analytics. However, dimensionality reduction can also be a useful method for data exploration, but it was only used in one study [57]. One reason for this could be the irregular temporal nature of AHR-data; while some methods have been developed to reduce the dimension of temporal data [107, 108, 109], we did not observe any applied to AHRs.

Can high-level healthcare concepts be modelled as latent variables?

The indirect nature of AHR data was a limitation recognised by some of the studies explored in this review. These studies seek to model the unobserved, latent variables present in AHRs using a variety of techniques, such as Hidden Markov Models (HMMs) [73, 74, 76] and topic models [56, 63, 85]. Future research should examine how these latent variables relate to various high-level concepts such as analysis of complex healthcare patterns or medical trajectories and pathways since these healthcare concepts are not recorded in AHRs.

How can AHR analytics insights be validated and interpreted?

Clinical data mining models must be validated to a standard that is much higher than in many other fields. Typically, data mining validation commonly used other fields (such as digital marketing or retail) seeks to demonstrate reproducible performance using standard techniques such as cross-validation. However, validity in healthcare must be applied at a cohort-level.

With recent advances in machine learning (e.g., deep learning), understanding a model's decision process is becoming increasingly difficult. Healthcare is a field where the interpretability of a model's decision process is critical. There are recent advancements in ExplainableAI, with recent studies such as [110] that use Shapley value based interpretations, proposing workarounds for black-box nature of state. Interpretability leads to trustworthiness [111] and constructing trustworthy models is advantageous for the adoption of machine-learned models in clinical practice. Novel methods are needed to validate and interpret ML models.

How can the interoperability of AHRs be maximised?

Health standards, vocabulary, and practices vary across different countries and continents. Achieving true data interoperability requires the development and implementation of standards and clinical-content models for the unambiguous representation and exchange of clinical meaning. Various

international certification and standards bodies pursue this goal and maintain resources such as international clinical terminology (CT)⁴ and International Disease Classification (ICD)⁵ [99]. AHR-based research should pursue solutions that are interoperable across these different standards. AHR researchers and data custodians should apply the FAIR data principles [112] in pursuit of this goal.

6. Discussion

This study highlights the trends in AHR-based research. Clustering, neural networks, and network analysis were among the most widely applied techniques used for analysing AHRs.

The most common application of AHR-based research is for understanding and predicting patient's medical trajectories. Other common applications include the construction of evidence-based guidelines, analysing patient comorbidities or cohorts and for biomarker discovery.

A common challenge for the identified ML methods is the use of abstract features, that make it difficult to explore the complex features within AHRs. Neural networks were the most widely applied technique to overcome this challenge, seeking to learn informative, encoded representations of AHRs. Neural networks struggle however in generating easily interpretable results, with learned representations requiring further downstream processing and analytics. Other studies address the abstract feature issue by modelling latent variables in AHRs with Markov-based techniques, topic modelling, or Bayesian networks. These studies typically use patient features, such as health state or treatment behaviours, to understand clinical pathways.

While this systematic review focuses specifically on AHRs, adjacent areas of research that use AHRs, such as the use of AHRs in combination with other EHR modalities or with genomic data, are not within the scope of this study. Using a diverse set of data modalities could contribute to a better understanding of health service patterns, but may also require novel analytics techniques and methods for combining knowledge from differing data modalities.

7. Conclusion

In this study, we identify the significance of AHRs in emerging health informatics research, and provide a considered definition for the term. AHRs are principally used to pursue health service pattern mining applications, particularly for learning about population-health.

This systematic review answers four key research questions concerning the use of AHRs for HSPM. We identify twelve machine learning techniques that are used to analyse AHRs, and eight key health informatics applications. Furthermore, we analyse how each technique is applied

⁴SNOMED clinical terminology from the International Health Terminology Standards Development Organisation (IHTSDO) <http://www.ihtsdo.org>

⁵ICD codes are standardised by the World Health Organisation (WHO) <http://www.who.int/classifications/icd/en>

and which techniques are used in pursuit of each application, revealing the current landscape of AHR-based research. Finally, this analysis highlighted many data-related and technique-related limitations of AHR-based research. We also identify six open research questions to provide clarity on the current state of AHR-based research and future research opportunities in this area.

It is clear that the use of AHRs in health informatics research is becoming more commonplace among researchers and healthcare institutions such as hospitals and governments. At the same time, machine learning techniques are advancing and are consequently being used to pursue increasingly complex research objectives. These trends are reaffirmed by the analysis in this review, with emerging machine learning techniques being used to analyse latent and complex healthcare patterns in AHRs.

Acknowledgements

This project was supported by a Commonwealth Services Contract (CA-DATA-200421) funded by Cancer Australia.

The authors thank Massimo Piccardi for providing feedback on a draft of this manuscript.

References

- [1] P. K. D. Pramanik, S. Pal, M. Mukhopadhyay, Healthcare big data: A comprehensive overview, *Intelligent systems for healthcare management and delivery* (2019) 72–100.
- [2] S. M. Shah, R. A. Khan, Secondary use of electronic health record: Opportunities and challenges, *IEEE Access* 8 (2020) 136947–136965.
- [3] P. Yadav, M. Steinbach, V. Kumar, G. Simon, Mining electronic health records (EHRs) a survey, *ACM Computing Surveys (CSUR)* 50 (2018) 1–40.
- [4] J. Chen, W. Wei, C. Guo, L. Tang, L. Sun, Textual analysis and visualization of research trends in data mining for electronic health records, *Health Policy and Technology* 6 (2017) 389–400.
- [5] A. P. Kurniati, O. Johnson, D. Hogg, G. Hall, Process mining in oncology: A literature review, in: 2016 6th International Conference on Information Communication and Management (ICICM), IEEE, 2016, pp. 291–297.
- [6] E. Rojas, J. Munoz-Gama, M. Sepúlveda, D. Capurro, Process mining in healthcare: A literature review, *Journal of Biomedical Informatics* 61 (2016) 224–236.
- [7] T. G. Erdogan, A. Tarhan, Systematic mapping of process mining studies in healthcare, *IEEE Access* 6 (2018) 24543–24567.
- [8] A. Guzzo, A. Rullo, E. Vocaturo, Process mining applications in the healthcare domain: A comprehensive review, *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 12 (2022) e1442.
- [9] J. Munoz-Gama, N. Martin, C. Fernandez-Llatas, O. A. Johnson, M. Sepúlveda, E. Helm, V. Galvez-Yanjari, E. Rojas, A. Martinez-Millana, D. Aloini, et al., Process mining for healthcare: Characteristics and challenges, *Journal of Biomedical Informatics* 127 (2022) 103994.
- [10] J. C. Brunson, R. C. Laubenbacher, Applications of network analysis to routinely collected health care data: a systematic review, *Journal of the American Medical Informatics Association* 25 (2017) 210–221. URL: <https://doi.org/10.1093/Fjamia%2Focx052>. doi:10.1093/famia/ocx052.
- [11] B. Shickel, P. J. Tighe, A. Bihorac, P. Rashidi, Deep EHR: a survey of recent advances in deep learning techniques for electronic health

- record (EHR) analysis, *IEEE Journal of Biomedical and Health Informatics* 22 (2017) 1589–1604.
- [12] C. Xiao, E. Choi, J. Sun, Opportunities and challenges in developing deep learning models using electronic health records data: a systematic review, *Journal of the American Medical Informatics Association* 25 (2018) 1419–1428.
- [13] S. M. Cadarette, L. Wong, An introduction to health care administrative data, *The Canadian Journal of Hospital Pharmacy* 68 (2015) 232.
- [14] Australian Institute of Health and Wellness (AIHW), Our data collections, <https://www.aihw.gov.au/about-our-data/our-data-collections>, 2021. Accessed: 2022-04-11.
- [15] National Health Service (NHS), List of administrative sources, <https://digital.nhs.uk/data-and-information/find-data-and-publications/statement-of-administrative-sources/list-of-administrative-sources>, 2021. Accessed: 2022-04-11.
- [16] D. Kindig, G. Stoddart, What is population health?, *American journal of public health* 93 (2003) 380–383.
- [17] Á. Rebuge, D. R. Ferreira, Business process analysis in healthcare environments: A methodology based on process mining, *Information systems* 37 (2012) 99–116.
- [18] M. Newman, *Networks*, Oxford University Press, 2018.
- [19] K. Petersen, R. Feldt, S. Mujtaba, M. Mattsson, Systematic mapping studies in software engineering, in: *12th International Conference on Evaluation and Assessment in Software Engineering (EASE)* 12, 2008, pp. 1–10.
- [20] J. D. Harris, C. E. Quatman, M. Manring, R. A. Siston, D. C. Flanigan, How to write a systematic review, *The American journal of sports medicine* 42 (2014) 2761–2768.
- [21] C. Wohlin, Guidelines for snowballing in systematic literature studies and a replication in software engineering, in: *Proceedings of the 18th international conference on evaluation and assessment in software engineering*, 2014, pp. 1–10.
- [22] M. J. Page, J. E. McKenzie, P. M. Bossuyt, I. Boutron, T. C. Hoffmann, C. D. Mulrow, L. Shamseer, J. M. Tetzlaff, E. A. Akl, S. E. Brennan, R. Chou, J. Glanville, J. M. Grimshaw, A. Hróbjartsson, M. M. Lalu, T. Li, E. W. Loder, E. Mayo-Wilson, S. McDonald, L. A. McGuinness, L. A. Stewart, J. Thomas, A. C. Tricco, V. A. Welch, P. Whiting, D. Moher, The PRISMA 2020 statement: an updated guideline for reporting systematic reviews, *British Medical Journal* 372 (2021). URL: <https://www.bmj.com/content/372/bmj.n71>. doi:10.1136/bmj.n71. arXiv:<https://www.bmj.com/content/372/bmj.n71.full.pdf>.
- [23] K. Khan, R. Kunz, J. Kleijnen, G. Antes, *Systematic reviews to support evidence-based medicine*, CRC press, 2011.
- [24] E. Choi, M. T. Bahadori, E. Searles, C. Coffey, M. Thompson, J. Bost, J. Tejedor-Sojo, J. Sun, Multi-layer representation learning for medical concepts, in: *proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1495–1504.
- [25] Z. Huang, W. Dong, P. Bath, L. Ji, H. Duan, On mining latent treatment patterns from electronic medical records, *Data mining and knowledge discovery* 29 (2015) 914–949.
- [26] A. B. Jensen, P. L. Moseley, T. I. Oprea, S. G. Ellesøe, R. Eriksson, H. Schmock, P. B. Jensen, L. J. Jensen, S. Brunak, Temporal disease trajectories condensed from population-wide registry data covering 6.2 million patients, *Nature communications* 5 (2014) 1–10.
- [27] C. Liu, F. Wang, J. Hu, H. Xiong, Temporal phenotyping from longitudinal electronic health records: A graph based framework, in: *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, 2015, pp. 705–714.
- [28] K. Baker, E. Dunwoodie, R. G. Jones, A. Newsham, O. Johnson, C. P. Price, J. Wolstenholme, J. Leal, P. McGinley, C. Twelves, et al., Process mining routinely collected electronic health records to define real-life clinical pathways during chemotherapy, *International journal of medical informatics* 103 (2017) 32–41.
- [29] E. Choi, M. T. Bahadori, L. Song, W. F. Stewart, J. Sun, GRAM: graph-based attention model for healthcare representation learning, in: *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, 2017, pp. 787–795.
- [30] E. Choi, C. Xiao, W. F. Stewart, J. Sun, MiME: multilevel medical embedding of electronic health records for predictive healthcare, in: *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018, pp. 4552–4562.
- [31] R. Guo, T. Fujiwara, Y. Li, K. M. Lima, S. Sen, N. K. Tran, K.-L. Ma, Comparative visual analytics for assessing medical records with sequence embedding, *Visual Informatics* 4 (2020) 72–85.
- [32] T. Pham, T. Tran, D. Phung, S. Venkatesh, Predicting healthcare trajectories from medical records: A deep learning approach, *Journal of biomedical informatics* 69 (2017) 218–229.
- [33] J. Zhang, K. Kowsari, J. H. Harrison, J. M. Lobo, L. E. Barnes, Patient2vec: A personalized interpretable deep representation of the longitudinal electronic health record, *IEEE Access* 6 (2018) 65333–65346.
- [34] B. Jin, H. Yang, L. Sun, C. Liu, Y. Qu, J. Tong, A treatment engine by predicting next-period prescriptions, in: *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 1608–1616.
- [35] Y. Li, Z. Zhu, H. Wu, S. Ding, Y. Zhao, CCAE: Cross-field categorical attributes embedding for cancer clinical endpoint prediction, *Artificial Intelligence in Medicine* 107 (2020) 101915.
- [36] X. Xu, Y. Wang, T. Jin, J. Wang, Learning the representation of medical features for clinical pathway analysis, in: *International Conference on Database Systems for Advanced Applications*, Springer, 2018, pp. 37–52.
- [37] B. K. Beaulieu-Jones, P. Orzechowski, J. H. Moore, Mapping patient trajectories using longitudinal extraction and deep learning in the MIMIC-III critical care database, in: *Pacific Symposium on Biocomputing 2018: Proceedings of the Pacific Symposium*, World Scientific, 2018, pp. 123–132.
- [38] S. Hong, M. Wu, H. Li, Z. Wu, Event2Vec: Learning representations of events on temporal sequences, in: *Web and Big Data*, Springer International Publishing, 2017, pp. 33–47. URL: https://doi.org/10.1007/978-3-319-63564-4_3. doi:10.1007/978-3-319-63564-4_3.
- [39] L. Li, R. Zuo, A. Coston, J. C. Weiss, G. H. Chen, Neural topic models with survival supervision: Jointly predicting time-to-event outcomes and learning how clinical features relate, in: *Artificial Intelligence in Medicine*, Springer International Publishing, 2020, pp. 371–381. URL: https://doi.org/10.1007/978-3-030-59137-3_33. doi:10.1007/978-3-030-59137-3_33.
- [40] X. Lu, L. Cui, Z. Sun, Y. Zhu, ProAID: path-based reasoning for self-attentional disease prediction, *Knowledge and Information Systems* 63 (2021) 3087–3101. URL: https://doi.org/10.1007/978-3-030-59137-3_33. doi:10.1007/s10115-021-01617-w.
- [41] J. Wolff, A. Gary, D. Jung, C. Normann, K. Kaier, H. Binder, K. Domschke, A. Klimke, M. Franz, Predicting patient outcomes in psychiatric hospitals with routine data: a machine learning approach, *BMC Medical Informatics and Decision Making* 20 (2020) 1–9. URL: <https://doi.org/10.1186/s12911-020-1042-2>. doi:10.1186/s12911-020-1042-2.
- [42] E. Steinberg, K. Jung, J. A. Fries, C. K. Corbin, S. R. Pfohl, N. H. Shah, Language models are an effective representation learning technique for electronic health record data, *Journal of Biomedical Informatics* 113 (2021) 103637.
- [43] H. Zheng, I. O. Ryzhov, W. Xie, J. Zhong, Personalized multimorbidity management for patients with type 2 diabetes using reinforcement learning of electronic health records, *Drugs* 81 (2021) 471–482. URL: <https://doi.org/10.1007/s40265-020-01435-4>. doi:10.1007/s40265-020-01435-4.
- [44] A. Caruana, M. Bandara, D. Catchpoole, P. J. Kennedy, Beyond topics: Discovering latent healthcare objectives from event sequences, in: G. Long, X. Yu, S. Wang (Eds.), *AI 2021: Advances in Artificial Intelligence*, Springer International Publishing, Cham, 2022, pp. 368–380.

- [45] L. Gerrard, X. Peng, A. Clarke, C. Schlegel, J. Jiang, Predicting outcomes for cancer patients with transformer-based multi-task learning, in: G. Long, X. Yu, S. Wang (Eds.), *AI 2021: Advances in Artificial Intelligence*, Springer International Publishing, Cham, 2022, pp. 381–392.
- [46] J. G. D. Ochoa, F. E. Mustafa, Graph neural network modelling as a potentially effective method for predicting and analyzing procedures based on patients' diagnoses, *Artificial Intelligence in Medicine* 131 (2022) 102359.
- [47] F. Doshi-Velez, Y. Ge, I. Kohane, Comorbidity clusters in autism spectrum disorders: An electronic health record time-series analysis, *Pediatrics* 133 (2014).
- [48] Y. Zhang, R. Padman, L. Wasserman, N. Patel, P. Teredesai, Q. Xie, On clinical pathway discovery from electronic health record data, *IEEE Intelligent Systems* 30 (2015) 70–75. URL: <https://doi.org/10.1109/2Fmis.2015.14>. doi:10.1109/mis.2015.14.
- [49] F. S. Roque, P. B. Jensen, H. Schmock, M. Dalgaard, M. Andreatta, T. Hansen, K. Søbey, S. Bredkjær, A. Juul, T. Werge, et al., Using electronic patient records to discover disease correlations and stratify patient cohorts, *PLOS Computational Biology* 7 (2011) e1002141.
- [50] C. Sideris, M. Pourhomayoun, H. Kalantarian, M. Sarrafzadeh, A flexible data-driven comorbidity feature extraction framework, *Computers in biology and medicine* 73 (2016) 165–172.
- [51] L. Chen, N. Blumm, N. Christakis, A. Barabasi, T. S. Deisboeck, Cancer metastasis networks and the prediction of progression patterns, *British journal of cancer* 101 (2009) 749–758.
- [52] J. Chen, L. Sun, C. Guo, W. Wei, Y. Xie, A data-driven framework of typical treatment process extraction and evaluation, *Journal of biomedical informatics* 83 (2018) 178–195.
- [53] A. C. Apunike, L. Oliveira-Ciabati, T. L. Sanches, L. L. de Oliveira, M. N. Sanchez, R. M. Galliez, D. Alves, Analyses of public health databases via clinical pathway modelling: TBWEB, in: *International Conference on Computational Science*, Springer, 2020, pp. 550–562.
- [54] H. Johns, J. Hearne, J. Bernhardt, L. Churilov, Clustering clinical and health care processes using a novel measure of dissimilarity for variable-length sequences of ordinal states, *Statistical methods in medical research* 29 (2020) 3059–3075.
- [55] R. J. C. Bose, W. M. van der Aalst, Trace clustering based on conserved patterns: Towards achieving better process models, in: *International Conference on Business Process Management*, Springer, 2009, pp. 170–181.
- [56] E. S. Prokofyeva, R. D. Zaytsev, S. V. Maltseva, Application of modern data analysis methods to cluster the clinical pathways in urban medical facilities, in: *2019 IEEE 21st Conference on Business Informatics (CBI)*, volume 1, IEEE, 2019, pp. 75–83.
- [57] J. Chen, C. Guo, L. Sun, M. Lu, Mining typical treatment duration patterns for rational drug use from electronic medical records, *Journal of Systems Science and Systems Engineering* 28 (2019) 602–620.
- [58] E. Aspland, P. R. Harper, D. Gartner, P. Webb, P. Barrett-Lee, Modified Needleman–Wunsch algorithm for clinical pathway clustering, *Journal of Biomedical Informatics* 115 (2021) 103668.
- [59] D. M. Bean, C. Stringer, N. Beeknoo, J. Teo, R. J. Dobson, Network analysis of patient flow in two uk acute care hospitals identifies key sub-networks for A&E performance, *PloS one* 12 (2017) e0185912.
- [60] B. Hompes, J. Buijs, W. Van der Aalst, P. Dixit, J. Buurman, Discovering deviating cases and process variants using trace clustering, in: *Proceedings of the 27th Benelux Conference on Artificial Intelligence (BNAIC)*, November, 2015, pp. 5–6.
- [61] M. Chambard, T. Guyet, Y.-L. NGuyen, E. Audureau, Temporal phenotyping for characterisation of hospital care pathways of COVID19 patients, in: *Advanced Analytics and Learning on Temporal Data*, Springer International Publishing, 2021, pp. 55–70. URL: https://doi.org/10.1007/2F978-3-030-91445-5_4. doi:10.1007/978-3-030-91445-5_4.
- [62] K. N. M. Kumar, S. Sampath, M. Imran, N. Pradeep, Clustering diagnostic codes: Exploratory machine learning approach for preventive care of chronic diseases, in: *Advances in Intelligent Systems and Computing*, Springer Singapore, 2020, pp. 551–564. URL: https://doi.org/10.1007/2F978-981-15-5679-1_53. doi:10.1007/978-981-15-5679-1_53.
- [63] C.-W. Huang, R. Lu, U. Iqbal, S.-H. Lin, P. A. Nguyen, H.-C. Yang, C.-F. Wang, J. Li, K.-L. Ma, Y.-C. Li, W.-S. Jian, A richly interactive exploratory data analysis and visualization tool using electronic medical records, *BMC Medical Informatics and Decision Making* 15 (2015). URL: <https://doi.org/10.1186/2Fs12911-015-0218-7>. doi:10.1186/s12911-015-0218-7.
- [64] K. Steinhäuser, N. V. Chawla, A network-based approach to understanding and predicting diseases, in: *Social computing and behavioral modeling*, Springer, 2009, pp. 1–8.
- [65] D. A. Hanauer, N. Ramakrishnan, Modeling temporal relationships in large scale clinical associations, *Journal of the American Medical Informatics Association* 20 (2013) 332–341.
- [66] B. S. Glicksberg, L. Li, M. A. Badgeley, K. Shameer, R. Kosoy, N. D. Beckmann, N. Pho, J. Hakenberg, M. Ma, K. L. Ayers, G. E. Hoffman, S. Dan Li, E. E. Schadt, C. J. Patel, R. Chen, J. T. Dudley, Comparative analyses of population-scale phenomic data in electronic medical records reveal race-specific disease networks, *Bioinformatics* 32 (2016) i101–i110. URL: <https://doi.org/10.1093/bioinformatics/btw282>. doi:10.1093/bioinformatics/btw282.
- [67] V. Kannan, F. Swartz, N. A. Kiani, G. Silberberg, G. Tspiras, D. Gomez-Cabrero, K. Alexanderson, J. Tegnèr, Conditional disease development extracted from longitudinal health care cohort data using layered network construction, *Scientific reports* 6 (2016) 1–14.
- [68] W. Dong, E. W. Lee, V. S. Hertzberg, R. L. Simpson, J. C. Ho, GASP: Graph-based approximate sequential pattern mining for electronic health records, in: *New Trends in Database and Information Systems*, Springer International Publishing, 2021, pp. 50–60. URL: https://doi.org/10.1007/2F978-3-030-85082-1_5. doi:10.1007/978-3-030-85082-1_5.
- [69] M. Kushima, Y. Honda, H. H. Le, T. Yamazaki, K. Araki, H. Yokota, Extraction and graph structuring of variants by detecting common parts of frequent clinical pathways, in: *International MultiConference of Engineers and Computer Scientists*, 2018, pp. 207–218.
- [70] S. Zhang, L. Liu, H. Li, L. Cui, Collaborative prediction model of disease risk by mining electronic health records, in: *Collaborate Computing: Networking, Applications and Worksharing*, Springer International Publishing, 2017, pp. 71–82. URL: https://doi.org/10.1007/2F978-3-319-59288-6_7. doi:10.1007/978-3-319-59288-6_7.
- [71] Y. fei Wang, J. jing Wang, W. Peng, Y. hao Ren, C. Gao, Y. lun Li, R. Wang, X. feng Wang, S. jun Han, J. yu Lyu, J. ming Huan, C. Chen, H. yan Wang, Z. xin Shu, X. zhong Zhou, W. Li, Identification of hypertension subgroups through topological analysis of symptom-based patient similarity, *Chinese Journal of Integrative Medicine* 27 (2021) 656–665. URL: <https://doi.org/10.1007/s11655-021-3336-3>. doi:10.1007/s11655-021-3336-3.
- [72] K. Maass, M. Kim, A markov decision process approach to optimizing cancer therapy using multiple modalities, *Mathematical Medicine and Biology: a journal of the IMA* 37 (2020) 22–39.
- [73] Z. Huang, Z. Ge, W. Dong, K. He, H. Duan, Probabilistic modeling personalized treatment pathways using electronic health records, *Journal of biomedical informatics* 86 (2018) 33–48.
- [74] A. Leontjeva, R. Conforti, C. Di Francescomarino, M. Dumas, F. M. Maggi, Complex symbolic sequence encodings for predictive monitoring of business processes, in: *International Conference on Business Process Management*, Springer, 2016, pp. 297–313.
- [75] S. Nagrecha, P. B. Thomas, K. Feldman, N. V. Chawla, Predicting chronic heart failure using diagnoses graphs, in: *Lecture Notes in Computer Science*, Springer International Publishing, 2017, pp. 295–312. URL: https://doi.org/10.1007/2F978-3-319-66808-6_20. doi:10.1007/978-3-319-66808-6_20.
- [76] M. L. P. Bueno, A. Hommersom, P. J. F. Lucas, M. Lobo, P. P. Rodrigues, Modeling the dynamics of multiple disease occurrence by latent states, in: *Lecture Notes in Computer Science*, Springer International Publishing, 2018, pp. 93–107. URL: <https://doi.org/>

- 10.1007/978-3-030-00461-3_7. doi:10.1007/978-3-030-00461-3_7.
- [77] D. Roder, G. W. Zhao, S. Challam, A. Little, E. Elder, G. Kostadinovska, L. Woodland, D. Currow, Female breast cancer in New South Wales, Australia, by country of birth: implications for health-service delivery, *BMC public health* 21 (2021) 1–14.
- [78] L. Te Marvelde, P. McNair, K. Whitfield, P. Autier, P. Boyle, R. Sullivan, R. J. Thomas, Alignment with indices of a care pathway is associated with improved survival: An observational population-based study in colon cancer patients, *EClinicalMedicine* 15 (2019) 42–50.
- [79] Z. Shahabi-Kargar, A. Johnston, M. Warner-Smith, N. Creighton, D. Roder, Differences in breast cancer treatment pathways for women participating in screening through BreastScreen New South Wales (BSNSW), *Australasian Medical Journal* 13 (2020).
- [80] C. Li, S. Gupta, S. Rana, W. Luo, S. Venkatesh, D. Ashely, D. Phung, Toxicity prediction in cancer using multiple instance learning in a multi-task framework, in: *Advances in Knowledge Discovery and Data Mining*, Springer International Publishing, 2016, pp. 152–164. URL: https://doi.org/10.1007/978-3-319-31753-3_13. doi:10.1007/978-3-319-31753-3_13.
- [81] X. Sun, A. Douiri, M. Gulliford, Applying machine learning algorithms to electronic health records to predict pneumonia after respiratory tract infection, *Journal of Clinical Epidemiology* 145 (2022) 154–163.
- [82] I. Kaur, M. Doja, T. Ahmad, Time-range based sequential mining for survival prediction in prostate cancer, *Journal of Biomedical Informatics* 110 (2020) 103550.
- [83] H. Estiri, S. Vasey, S. N. Murphy, Transitive sequential pattern mining for discrete clinical data, in: *Artificial Intelligence in Medicine*, Springer International Publishing, 2020, pp. 414–424. URL: https://doi.org/10.1007/978-3-030-59137-3_37. doi:10.1007/978-3-030-59137-3_37.
- [84] A. Vincent-Paulraj, G. Burnside, F. Coenen, M. Pirmohamed, L. Walker, Sequential association rule mining revisited: A study directed at relational pattern mining for multi-morbidity, in: *Lecture Notes in Computer Science*, Springer International Publishing, 2021, pp. 241–253. URL: https://doi.org/10.1007/978-3-030-91100-3_20. doi:10.1007/978-3-030-91100-3_20.
- [85] Z. Huang, W. Dong, L. Ji, H. Duan, Outcome prediction in clinical treatment processes, *Journal of Medical Systems* 40 (2015). URL: https://doi.org/10.1007/978-3-030-91100-3_20. doi:10.1007/s10916-015-0380-6.
- [86] M. R. Boland, N. P. Tatonetti, G. Hripcsak, Development and validation of a classification approach for extracting severity automatically from electronic health records, *Journal of Biomedical Semantics* 6 (2015). URL: <https://doi.org/10.1186/2Fs13326-015-0010-8>. doi:10.1186/s13326-015-0010-8.
- [87] Y. Maali, O. Perez-Concha, E. Coiera, D. Roffe, R. O. Day, B. Gallego, Predicting 7-day, 30-day and 60-day all-cause unplanned readmission: a case study of a Sydney hospital, *BMC Medical Informatics and Decision Making* 18 (2018). URL: <https://doi.org/10.1186/2Fs12911-017-0580-8>. doi:10.1186/s12911-017-0580-8.
- [88] K.-J. Wang, J.-L. Chen, K.-H. Chen, K.-M. Wang, Survivability prognosis for lung cancer patients at different severity stages by a risk factor-based Bayesian network modeling, *Journal of Medical Systems* 44 (2020) 65.
- [89] J. C. Weiss, D. Page, Forest-based point process for event prediction from electronic health records, in: *Advanced Information Systems Engineering*, Springer Berlin Heidelberg, 2013, pp. 547–562. URL: https://doi.org/10.1007/978-3-642-40994-3_35. doi:10.1007/978-3-642-40994-3_35.
- [90] G. Du, Y. Shi, A. Liu, T. Liu, Variance risk identification and treatment of clinical pathway by integrated Bayesian network and association rules mining, *Entropy* 21 (2019) 1191.
- [91] F. Wang, N. Lee, J. Hu, J. Sun, S. Ebadollahi, A. F. Laine, A framework for mining signatures from event sequences and its applications in healthcare data, *IEEE transactions on pattern analysis and machine intelligence* 35 (2012) 272–285.
- [92] D. Nguyen, W. Luo, D. Phung, S. Venkatesh, Understanding toxicities and complications of cancer treatment: A data mining approach, in: *AI 2015: Advances in Artificial Intelligence*, Springer International Publishing, 2015, pp. 431–443. URL: https://doi.org/10.1007/978-3-319-26350-2_38. doi:10.1007/978-3-319-26350-2_38.
- [93] J. Chen, L. Sun, C. Guo, Y. Xie, A fusion framework to extract typical treatment patterns from electronic medical records, *Artificial intelligence in medicine* 103 (2020) 101782.
- [94] L. J. Liu, H. Zhang, J. Di, J. Chen, ELMV: an ensemble-learning approach for analyzing electrical health records with significant missing values, in: *Proceedings of the 11th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics*, 2020, pp. 1–10.
- [95] Y. Xue, D. Klabjan, Y. Luo, Mixture-based multiple imputation model for clinical data with a temporal dimension, in: *2019 IEEE International Conference on Big Data (Big Data)*, IEEE, 2019, pp. 245–252.
- [96] S. Huda, J. Yearwood, H. F. Jelinek, M. M. Hassan, G. Fortino, M. Buckland, A hybrid feature selection with ensemble classification for imbalanced healthcare data: A case study for brain tumor diagnosis, *IEEE access* 4 (2016) 9145–9154.
- [97] P. Ray, J. Wimalasiri, The need for technical solutions for maintaining the privacy of EHR, in: *2006 International Conference of the IEEE Engineering in Medicine and Biology Society*, 2006, pp. 4686–4689. doi:10.1109/IEMBS.2006.260862.
- [98] L. Myers, J. Stevens, Using EHR to Conduct Outcome and Health Services Research, Springer International Publishing, Cham, 2016, pp. 61–70. URL: https://doi.org/10.1007/978-3-319-43742-2_7. doi:10.1007/978-3-319-43742-2_7.
- [99] World Health Organization, Icd-10 : international statistical classification of diseases and related health problems : tenth revision, 2004.
- [100] Healthcare Cost and Utilization Project, Clinical Classifications Software (CCS) for ICD-10-PCS (beta version), <https://www.hcup-us.ahrq.gov/toolssoftware/ccs10/ccs10.jsp>, 2019. Accessed: 2022-03-11.
- [101] A. E. Johnson, T. J. Pollard, L. Shen, H. L. Li-Wei, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. A. Celi, R. G. Mark, MIMIC-III, a freely accessible critical care database, *Scientific data* 3 (2016) 1–9.
- [102] A. Johnson, L. Bulgarelli, T. Pollard, S. Horng, L. A. Celi, R. Mark IV, MIMIC-IV (version 0.4), *PhysioNet* (2020).
- [103] B. van Dongen, Real-life event logs - hospital log, 4TU.ResearchData.Dataset (2011). doi:10.4121/uuid:d9769f3d-0ab0-4fb8-803b-0d1120ffc54.
- [104] T. E. Raghunathan, Synthetic data, *Annual Review of Statistics and Its Application* 8 (2021) 129–140. URL: <https://doi.org/10.1146/annurev-statistics-040720-031848>. doi:10.1146/annurev-statistics-040720-031848.
- [105] A. Goncalves, P. Ray, B. Soper, J. Stevens, L. Coyle, A. P. Sales, Generation and evaluation of synthetic patient data, *BMC Medical Research Methodology* 20 (2020). URL: <https://doi.org/10.1186/2Fs12874-020-00977-1>. doi:10.1186/s12874-020-00977-1.
- [106] A. Tucker, Z. Wang, Y. Rotalinti, P. Myles, Generating high-fidelity synthetic patient data for assessing machine learning healthcare software, *NPJ Digital Medicine* 3 (2020). URL: <https://doi.org/10.1038/2Fs41746-020-00353-9>. doi:10.1038/s41746-020-00353-9.
- [107] M. Gashler, T. Martinez, Temporal nonlinear dimensionality reduction, in: *The 2011 International Joint Conference on Neural Networks*, IEEE, 2011, pp. 1959–1966.
- [108] M. Ali, M. W. Jones, X. Xie, M. Williams, Timecluster: dimension reduction applied to temporal data for visual analytics, *The Visual Computer* 35 (2019) 1013–1026.
- [109] M. Lewandowski, J. Martinez-del Rincon, D. Makris, J.-C. Nebel, Temporal extension of Laplacian eigenmaps for unsupervised dimensionality reduction of time series, in: *2010 20th International Conference on Pattern Recognition*, IEEE, 2010, pp. 161–164.
- [110] Y. Liu, S. Qin, An interpretable machine learning approach for predicting hospital length of stay and readmission, in: *International*

Conference on Advanced Data Mining and Applications, Springer, 2022, pp. 73–85.

- [111] Z. C. Lipton, The myths of model interpretability: In machine learning, the concept of interpretability is both important and slippery., *Queue* 16 (2018) 31–57.
- [112] M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L. B. da Silva Santos, P. E. Bourne, et al., The fair guiding principles for scientific data management and stewardship, *Scientific data* 3 (2016) 1–9.