

FACULDADE DE MEDICINA DA UNIVERSIDADE DO PORTO

# **Knowledge Discovery in Healthcare: Exploring the role of real-world data to leverage clinical practice**

**João Filipe Coutinho de Almeida**

PREPARAÇÃO DA DISSERTAÇÃO



Programa Doutoral em Ciência de dados de saúde

Supervisor: Pedro Pereira Rodrigues

Second Supervisor: Ricardo Correia

November 3, 2023



Cover page



## Integrity Declaration



## Reproducibility

The code for all the experiments done in this thesis is stored online on GitHub. The link is as follows: <https://github.com/joofio/heads-thesis> From here, it should be possible to access the list of all the repositories involved in this thesis. Data is also available when possible. However, since most of the data used was directly retrieved from Electronic health record, it is blocked by ethical committees from sharing with third parties.





To my and my



# List of Publications

## Core Research Papers

The 7 papers described below are the core structure of this thesis (7 were already published, 1 is under review, and 2 are waiting for submission). The manuscripts are listed by order of appearance in the thesis.

Coutinho-Almeida, J., Cruz-Correia, R., & Rodrigues, P. (2022). Dataset Comparison Tool: Utility and Privacy. *Stud Health Technol Inform*, 294, 23–27.

Coutinho-Almeida, J., Rodrigues, P., & Cruz-Correia, R. (2021). GANs for Tabular Healthcare Data Generation: A Review on Utility and Privacy. In *Discovery Science* (pp. 282–291). Springer International Publishing.

## Other Publications and activities

In addition, during the duration of this thesis conduction, the candidate was also the author and co-author of other papers. Although these studies were not part of the thesis core structure, they were important to improve the researcher’s knowledge of the field and/or to present the results to the community. They are listed below:

Coutinho-Almeida, J., & Cruz-Correia, R. (2022). Developing a Process Mining Tool Based on HL7. *Procedia Computer Science*, 196, 501–508.

Holmgren, A., Esdar, M., Hüsters, J., & Coutinho-Almeida, J. (2023). Health Information Exchange: Understanding the Policy Landscape and Future of Data Interoperability. *Yearbook of Medical Informatics*, s-0043-1768719.

Costa, P., Almeida, J., Araujo, S., Alves, P., Cruz-Correia, R., Saranto, K., & Mantas, J. (2023). Biomedical and Health Informatics Teaching in Portugal: Current Status. *Heliyon*, 9(3).



# Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed vehicula lorem commodo dui. Fusce mollis feugiat elit. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Donec eu quam. Aenean consectetur odio quis nisi. Fusce molestie metus sed neque. Praesent nulla. Donec quis urna. Pellentesque hendrerit vulputate nunc. Donec id eros et leo ullamcorper placerat. Curabitur aliquam tellus et diam.

Ut tortor. Morbi eget elit. Maecenas nec risus. Sed ultricies. Sed scelerisque libero faucibus sem. Nullam molestie leo quis tellus. Donec ipsum. Nulla lobortis purus pharetra turpis. Nulla laoreet, arcu nec hendrerit vulputate, tortor elit eleifend turpis, et aliquam leo metus in dolor. Praesent sed nulla. Mauris ac augue. Cras ac orci. Etiam sed urna eget nulla sodales venenatis. Donec faucibus ante eget dui. Nam magna. Suspendisse sollicitudin est et mi.

Fusce sed ipsum vel velit imperdiet dictum. Sed nisi purus, dapibus ut, iaculis ac, placerat id, purus. Integer aliquet elementum libero. Phasellus facilisis leo eget elit. Nullam nisi magna, ornare at, aliquet et, porta id, odio. Sed volutpat tellus consectetur ligula. Phasellus turpis augue, malesuada et, placerat fringilla, ornare nec, eros. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Vivamus ornare quam nec sem mattis vulputate. Nullam porta, diam nec porta mollis, orci leo condimentum sapien, quis venenatis mi dolor a metus. Nullam mollis. Aenean metus massa, pellentesque sit amet, sagittis eget, tincidunt in, arcu. Vestibulum porta laoreet tortor. Nullam mollis elit nec justo. In nulla ligula, pellentesque sit amet, consequat sed, faucibus id, velit. Fusce purus. Quisque sagittis urna at quam. Ut eu lacus. Maecenas tortor nibh, ultricies nec, vestibulum varius, egestas id, sapien.

Phasellus ullamcorper justo id risus. Nunc in leo. Mauris auctor lectus vitae est lacinia egestas. Nulla faucibus erat sit amet lectus varius semper. Praesent ultrices vehicula orci. Nam at metus. Aenean eget lorem nec purus feugiat molestie. Phasellus fringilla nulla ac risus. Aliquam elementum aliquam velit. Aenean nunc odio, lobortis id, dictum et, rutrum ac, ipsum.

Ut tortor. Morbi eget elit. Maecenas nec risus. Sed ultricies. Sed scelerisque libero faucibus sem. Nullam molestie leo quis tellus. Donec ipsum. Nulla lobortis purus pharetra turpis. Nulla laoreet, arcu nec hendrerit vulputate, tortor elit eleifend turpis, et aliquam leo metus in dolor. Praesent sed nulla. Mauris ac augue. Cras ac orci. Etiam sed urna eget nulla sodales venenatis. Donec faucibus ante eget dui. Nam magna. Suspendisse sollicitudin est et mi.

Phasellus ullamcorper justo id risus. Nunc in leo. Mauris auctor lectus vitae est lacinia egestas. Nulla faucibus erat sit amet lectus varius semper. Praesent ultrices vehicula orci.

Ut tortor. Morbi eget elit. Maecenas nec risus. Sed ultricies. Sed scelerisque libero faucibus sem. Nullam molestie leo quis tellus. Donec ipsum.

**Keywords:** keyword1, Keyword2, keyword3



# Resumo

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed vehicula lorem commodo dui. Fusce mollis feugiat elit. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Donec eu quam. Aenean consectetur odio quis nisi. Fusce molestie metus sed neque. Praesent nulla. Donec quis urna. Pellentesque hendrerit vulputate nunc. Donec id eros et leo ullamcorper placerat. Curabitur aliquam tellus et diam.

Ut tortor. Morbi eget elit. Maecenas nec risus. Sed ultricies. Sed scelerisque libero faucibus sem. Nullam molestie leo quis tellus. Donec ipsum. Nulla lobortis purus pharetra turpis. Nulla laoreet, arcu nec hendrerit vulputate, tortor elit eleifend turpis, et aliquam leo metus in dolor. Praesent sed nulla. Mauris ac augue. Cras ac orci. Etiam sed urna eget nulla sodales venenatis. Donec faucibus ante eget dui. Nam magna. Suspendisse sollicitudin est et mi.

Fusce sed ipsum vel velit imperdiet dictum. Sed nisi purus, dapibus ut, iaculis ac, placerat id, purus. Integer aliquet elementum libero. Phasellus facilisis leo eget elit. Nullam nisi magna, ornare at, aliquet et, porta id, odio. Sed volutpat tellus consectetur ligula. Phasellus turpis augue, malesuada et, placerat fringilla, ornare nec, eros. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Vivamus ornare quam nec sem mattis vulputate. Nullam porta, diam nec porta mollis, orci leo condimentum sapien, quis venenatis mi dolor a metus. Nullam mollis. Aenean metus massa, pellentesque sit amet, sagittis eget, tincidunt in, arcu. Vestibulum porta laoreet tortor. Nullam mollis elit nec justo. In nulla ligula, pellentesque sit amet, consequat sed, faucibus id, velit. Fusce purus. Quisque sagittis urna at quam. Ut eu lacus. Maecenas tortor nibh, ultricies nec, vestibulum varius, egestas id, sapien.

Phasellus ullamcorper justo id risus. Nunc in leo. Mauris auctor lectus vitae est lacinia egestas. Nulla faucibus erat sit amet lectus varius semper. Praesent ultrices vehicula orci. Nam at metus. Aenean eget lorem nec purus feugiat molestie. Phasellus fringilla nulla ac risus. Aliquam elementum aliquam velit. Aenean nunc odio, lobortis id, dictum et, rutrum ac, ipsum.

Ut tortor. Morbi eget elit. Maecenas nec risus. Sed ultricies. Sed scelerisque libero faucibus sem. Nullam molestie leo quis tellus. Donec ipsum. Nulla lobortis purus pharetra turpis. Nulla laoreet, arcu nec hendrerit vulputate, tortor elit eleifend turpis, et aliquam leo metus in dolor. Praesent sed nulla. Mauris ac augue. Cras ac orci. Etiam sed urna eget nulla sodales venenatis. Donec faucibus ante eget dui. Nam magna. Suspendisse sollicitudin est et mi.

Phasellus ullamcorper justo id risus. Nunc in leo. Mauris auctor lectus vitae est lacinia egestas. Nulla faucibus erat sit amet lectus varius semper. Praesent ultrices vehicula orci.

Ut tortor. Morbi eget elit. Maecenas nec risus. Sed ultricies. Sed scelerisque libero faucibus sem. Nullam molestie leo quis tellus. Donec ipsum.

**Keywords:** keyword1, Keyword2, keyword3





# Acknowledgements

Queria aqui agradecer, fachabor, fachabor.

Author's Name



*“If you ain’t aim too high,  
Then you aim too low.”*

Jermaine Lamarr Cole



# Outline

The idea for this thesis first formed in my mind during a mental process of understanding how clinical knowledge could be improved in terms of quality, quantity, and speed of generation. The feeling was that new technology, especially the ones related to digital and informatics domain took years to be fully implemented in practice and harness the potential benefits they provided. I felt that healthcare, like other domains, had a serious gap between academia and industry. So the potential of all these discoveries was lost in "translation". So how could we leverage this?

This thesis is organized as follows:

Chapter 1 synthesizes the aim and specific objectives of this thesis. Chapter 2 presents a brief introduction to core concepts for the thesis, like Knowledge Discovery in Databases (KDD), Evidence Based Medicine (EBM) or privacy and ethical concerns.

Chapter 3 corresponds to the papers published. The papers cover a wide range of the traditional KDD steps so they are grouped around the phases they represent the most.

Chapter 4 represents the overall discussion of all of the papers and experiments done in the thesis.

Chapter 5 communicates the conclusion, limitations and future work.

Attachments include ethical permissions and supplementary data to some of the papers.



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Rationale . . . . .	1
1.2	Research Objectives . . . . .	2
<b>2</b>	<b>State of the art</b>	<b>3</b>
2.1	Artificial Intelligence . . . . .	3
2.2	Evidence Based Medicine . . . . .	4
2.3	Extracting Knowledge of Data . . . . .	6
2.4	Health Data Science . . . . .	8
2.5	Explainable Artificial Intelligence . . . . .	10
2.6	Causality . . . . .	13
2.7	Legal and ethical considerations . . . . .	17
<b>3</b>	<b>Case Studies</b>	<b>21</b>
3.1	Can GANs help create realistic datasets? . . . . .	21
3.1.1	Introduction . . . . .	21
3.1.2	Theoretical background . . . . .	22
3.1.3	Methods . . . . .	24
3.1.4	Results . . . . .	24
3.1.5	Implications for future research . . . . .	28
3.1.6	Conclusion . . . . .	28
3.2	Pulling the current metrics of assessing datasets . . . . .	28
3.2.1	Introduction . . . . .	29
3.2.2	Methods . . . . .	30
3.2.3	Results . . . . .	32
3.2.4	Discussion & Conclusion . . . . .	33
3.3	Can we use machine learning feature to compare datasets? . . . . .	33
3.3.1	Introduction . . . . .	33
3.3.2	Rationale and Related Work . . . . .	34
3.3.3	Materials & Methods . . . . .	35
3.3.3.1	Method Overview . . . . .	35
3.3.3.2	Data used . . . . .	37
3.3.4	Results . . . . .	37
3.3.5	Discussion . . . . .	41
3.3.6	Conclusion . . . . .	41
3.4	Data quality Metrics . . . . .	42
3.4.1	Introduction . . . . .	42
3.4.2	Background and Related Work . . . . .	43
3.4.3	Materials . . . . .	44
3.4.4	Methods . . . . .	44

3.4.5	Results . . . . .	46
3.4.6	Deployment & Validation . . . . .	48
3.4.7	Discussion . . . . .	48
3.4.8	Conclusion . . . . .	50
3.5	Leveraging Distributed systems in healthcare: is it advisable? . . . . .	51
3.5.1	Introduction . . . . .	51
3.5.2	Theoretical background and Related Work . . . . .	52
3.5.3	Materials . . . . .	53
3.5.4	Methods . . . . .	54
3.5.4.1	Model Performance Evaluation . . . . .	56
3.5.4.2	Model Training . . . . .	56
3.5.5	Results . . . . .	58
3.5.6	Discussion . . . . .	58
3.5.7	Conclusion . . . . .	63
3.6	Can Institutions share their performance metrics without hesitation of retaliation? . . . . .	63
3.6.1	Introduction . . . . .	64
3.6.2	Rationale and Related Work . . . . .	65
3.6.3	Materials & Methods . . . . .	66
3.6.3.1	Method Overview . . . . .	66
3.6.3.2	Data used . . . . .	67
3.6.4	Results . . . . .	67
3.6.5	Discussion . . . . .	68
3.6.6	Conclusion . . . . .	70
3.7	Leveraging data to assess treatment efficacy . . . . .	70
3.7.1	Introduction . . . . .	71
3.7.2	Materials & Methods . . . . .	72
3.7.3	Study Design . . . . .	72
3.7.4	Data collection . . . . .	72
3.7.5	Statistical Analysis . . . . .	73
3.7.6	Results . . . . .	74
3.7.7	Discussion . . . . .	76
3.7.8	Conclusion . . . . .	77
3.8	Leveraging data to create Clinical Decision Support Systems . . . . .	78
3.8.1	Introduction . . . . .	78
3.8.2	Rationale and Related Work . . . . .	79
3.8.3	Materials . . . . .	80
3.8.4	Methods . . . . .	80
3.8.5	Results . . . . .	82
3.8.5.1	The model . . . . .	82
3.8.5.2	Deployment . . . . .	83
3.8.5.3	Clinical Evaluation . . . . .	84
3.8.5.4	Potential Financial Impact . . . . .	84
3.8.6	Discussion . . . . .	86
3.8.7	Conclusion . . . . .	86
<b>4</b>	<b>Discussion</b>	<b>89</b>
<b>5</b>	<b>Limitations, future work and conclusions</b>	<b>93</b>



<b>A</b>	<b>95</b>
A.1 Data Dictionary . . . . .	96
<b>B</b>	<b>97</b>
B.1 C-section assessment questionnaire . . . . .	98
B.2 Data quality questionnaire . . . . .	99
<b>References</b>	<b>101</b>



## List of Figures

2.1	EBM adapted from [79] . . . . .	5
2.2	KDD Process, adapted from [69] . . . . .	7
3.1	Generative Adversarial Network (GAN) framework . . . . .	23
3.3	Continuous Variables plotted . . . . .	32
3.2	Categorical Variables plotted . . . . .	32
3.4	Cross-Validation of datasets . . . . .	35
3.5	Plot showing the decrease of the metric over increasingly changed datasets. . . .	38
3.6	Plot showing the values at 50% columns mutated across all datasets and algorithms per metric type . . . . .	39
3.7	Plot the variance of different repetitions for every metric and the number of differ- ent columns changed. X is the number of columns mutated. Colour is the number of repetitions for each mutation and Y is the variance of the data. . . . .	40
3.8	Result on a synthetic and real data . . . . .	41
3.9	Dimensions of data quality . . . . .	44
3.10	Network learned . . . . .	47
3.11	Workflow for creating the final score and which elements are used to do so. . . .	47
3.12	Model score for newly seen data . . . . .	49
3.13	Comparasion of clinical assessment of records with the model . . . . .	49
3.14	Heatmap of classification algorithm and silo vs Target variable and model type. Value is the AUROC mean of all 10 experiments. Y axis is the algorithm and silo. X axis is Target variable and Method. . . . .	59
3.15	Heatmap of regression algorithm and silo vs Target variable and model type. Value is the MAE mean of all 10 experiments. The y axis is the algorithm and silo. X axis is Target variable and Method. . . . .	60
3.16	Clustering for 3 continuous variables with 3 silos . . . . .	67
3.17	Clustering for 3 variables with 9 silos . . . . .	68
3.18	Clustering for 3 variables with 3 silos - (A) categorical variables with proportion with K-Means and (B) Categorical with K-modes . . . . .	70
3.19	Survival curves for Palbociclib and Ribociclib (1st line) - Progression Free Sur- vival and Overall Survival . . . . .	74
3.20	Survival curves (OS and PFS) comparing endocrine therapy (ET) to CDK4/6 in- hibitors as 1st line. p values shown as pairwise vs. ET. . . . .	76
3.21	Comparison of palbociclib and ribociclib survival curves adjusted for propensity scores . . . . .	77
3.22	Deployment and decision mechanism of the model . . . . .	83
3.23	Obstetrics questionnaires data . . . . .	85



## List of Tables

3.1	Summary of the articles selected. . . . .	25
3.2	Metrics utilised for evaluation . . . . .	26
3.3	Metrics Assessed . . . . .	31
3.4	Implemented Methods . . . . .	45
3.5	Validation Results: Column acronym with AUROC along with 95% CI . . . . .	46
3.6	Silos overview. . . . .	54
3.7	Silos overview part 2. . . . .	55
3.8	Metrics for centralised model, distributed model and local model . . . . .	58
3.9	Hypothesis testing of Distributed versus Centralised and local models . . . . .	61
3.10	Final Data points after convergence of clustering . . . . .	68
3.11	Final Data points after convergence and true centroids of the true means of each silo (TC) . . . . .	69
3.12	Descriptive statistics of cyclin-dependent kinase inhibitors group and endocrine therapy group. The Drug/combination refers to the actual drug or the combination for CDK4/6 . . . . .	73
3.13	Cox Regression with palbociclib and Ribociclib - Progression Free Survival and Overall Survival . . . . .	75
3.14	Distribution of feature used for prediction . . . . .	81
3.15	Distribution of Delivery Methods . . . . .	82
3.16	Performance Metrics in the training set . . . . .	82
3.17	Performance Metrics in the test set with chosen threshold . . . . .	82
3.18	Ruleset for financial support indexed to C-Sections. . . . .	85



# Abbreviations

**AI** Artificial Intelligence

**API** Application Programming Interface

**ATC** Anatomical Therapeutic Chemical

**ATE** Average Treatment Effect

**ATT** Average Treatment Effect on the Treated

**AUPRC** Area Under the Precision Recall Curve

**AUROC** Area Under the Receiver Operating Characteristic Curve

**BMI** Body Mass Index

**BN** Bayesian Network

**C-Section** Cesarean Section

**CausalML** Causal Machine Learning

**CDSS** Clinical Decision Support System

**CRISP-DM** Cross-Industry Standard Process for Data Mining

**DAG** Directed Acyclic Graph

**DPO** Data Protection Officer

**EBM** Evidence Based Medicine

**EDA** Exploratory Data Analysis

**EHDS** European Health Data Space

**EHR** Electronic health record

**EU** European Union

**FHIR** Fast Healthcare Interoperability Resources

**GAN** Generative Adversarial Network

**GDPR** General Data Protection Regulation

**HIPAA** Health Insurance Portability and Accountability Act

**HIS** Health Information System

<b>HL7</b>	Health Level Seven
<b>IKNL</b>	Integraal Kankercentrum Nederland
<b>IPTW</b>	Inverse Probability of Treatment Weighting
<b>IV</b>	Instrumental variable
<b>KDD</b>	Knowledge Discovery in Databases
<b>KNN</b>	K-Nearest Neighbours
<b>KS</b>	<i>Kolmogorov-Smirnov</i>
<b>LIME</b>	Local Interpretable Model-Agnostic Explanation
<b>MAE</b>	Mean Absolute Error
<b>MHRA</b>	Healthcare Products Regulatory Agency
<b>ML</b>	Machine Learning
<b>MRE</b>	Mean Relative Error
<b>PCA</b>	Principal Component Analysis
<b>POF</b>	Potential Outcome Framework
<b>RCT</b>	Randomized Clinical Trial
<b>RI</b>	Rand Index
<b>RMSE</b>	Root Mean Squared Error
<b>RWD</b>	Real World Data
<b>RWE</b>	Real World Evidence
<b>SCM</b>	Structural Causal Model
<b>SEM</b>	structural equation model
<b>SEMMA</b>	Sample, Explore, Modify, Model, and Assess
<b>SHAP</b>	SHapley Additive exPlanations
<b>SMOTE</b>	Synthetic Minority Oversampling Technique
<b>SVM</b>	Support Vector Machines
<b>UCI</b>	UC Irvine Machine Learning Repository
<b>USA</b>	United States of America
<b>VAE</b>	Variational Autoencoder
<b>XAI</b>	Explainable AI



## Glossary



*If we knew what it was we were doing, it would not be called research, would it?*

Albert Einstein

# 1

## Introduction

### 1.1 Rationale

Healthcare practice revolves a lot around technology. Technology in the definitional sense of referring to *"methods, systems, and devices which are the result of scientific knowledge being used for practical purposes"*. Healthcare and medicine are in fact applied sciences where we use our knowledge of biology, physics, chemistry, and math and apply those concepts in order to create treatments, diagnoses, procedures, etc. However, in the last 20-30 years, computer science and informatics started gaining traction in the healthcare space [8]. A paper-based industry is now being digitalized and computerized. This has been leading to an increase in the amount of data generated by healthcare systems [108, 140]. This data has the potential to greatly improve the current methods and practices in healthcare. However, it is still not being used to its full potential [108, 62]. This is especially important when we note that the gold standard of evidence creation is Randomized Clinical Trials (RCTs) which can vary on quality, time, and resources. A RCT may cost no less than 20 million euros to run, and according to a report submitted to the United States of America (USA) Department of Health and Human Services [174] can cost as much as 100 million USA dollars. This is indeed a very steep price to get the information we need to innovate. Parallel to this, usually supported by these RCTs are systematic reviews and meta-analyses, highly supported and promoted by Evidence Based Medicine (EBM) which are estimated to cost around 140 thousand dollars each [125]. Additionally, we have to take into account the time that it takes to create and publish a good paper on evidence synthesis, often making it hard to keep up with the pace of innovation.

So, we are now being faced with huge amounts of clinical data generated by Electronic health

records (EHRs) and Health Information Systems (HISs). But which tools are the most suited for harvesting the potential of this data? The capabilities and assumptions behind modern Knowledge Discovery in Databases (KDD), Machine Learning (ML), and Artificial Intelligence (AI) seem like a good approach for harvesting this potential. However, they are very different from the traditional statistical methods that are usually used in healthcare. So, in order to properly use these methods in healthcare and actually provide value to the patients, we need to understand the differences between these methods and how they can be used to complement each other.

Currently, we already have an idea of what are the major key areas that hinder the adoption of AI in healthcare like problems related to data privacy and security, data quality and integrity, interoperability, ethical considerations, and the fact that the hype of AI is far greater than the AI science, the acceptance, and trust of healthcare practitioners of AI based systems [129, 106], and how to properly evaluate the potential risks of AI in healthcare, just to mention a few [186]. This is a very complex problem that requires a lot of different approaches and solutions. It is a popular assumption that 87% of data science projects never get into production [4]. Even if numbers for the healthcare domain are not available at this time, it is safe to assume that the number is not much different, if not higher. And those that actually do, may never actually create any impact due to the lack of adoption by the healthcare practitioners or the lack of trust in the system [204].

So, with this introduction, we have there is still a long way to go to harvest all the potential healthcare data has to offer. And so my research objectives are focused on powering up this adoption. What can be done to improve these chances? What can we bring to the table to enhance the rate of success?

## 1.2 Research Objectives

This thesis has three main goals:

- Goal 1: Research methods for improving Data Quality. Whether through synthetic data generation to enlarge data volume and protect privacy (sections 3.1, 3.2 and 3.3) or by creating automatic data quality assessments (section 3.4)
- Goal 2: Assess alternative ways of usage of data without having access to all of it. This will be covered in sections 3.5 and 3.6.
- Goal 3: Difficulties and steps resulting from attempts to convert data into decisions and policies, whether through ML or traditional statistics. This will be covered in sections 3.7 and 3.8.

*The most exciting phrase to hear in science, the one that heralds new discoveries, is not 'Eureka!' but 'That's funny...'*

Isaac Asimov

# 2

## State of the art

### 2.1 Artificial Intelligence

AI has already been under public focus for a few years now, but its concept is still elusive, mainly due to the fact that the definition has been changing rapidly as well. From the very beginning, the field of AI was about not only understanding but also building intelligent entities [170]. Intelligent entities can be understood as machines that can act according to what is expected in a wide range of situations. The first work of AI could be credited to Warren McCulloch and Walter Pitts (1943) with the proposed model of artificial neurons. In the 50s, AI could be associated with the works of Christopher Strachey, of two chess-playing programs. In the 60s, the perceptrons could be indicated as state-of-the-art AI. In the 80s, expert systems were providing advanced reasoning that the so-called weak methods of previous iterations could not compete with. The 90s brought the probabilistic reasoning and ML which led to more robust systems that went further than the boolean logic used so far. In the 2000s, big data and ML got focused on. Big data was used as a matter symbolizing the increasing amounts of data in some industries [57], and ML as *the study of computer algorithms that improve automatically through experience* [127]. This last definition is especially important since it is currently used as a synonym of AI across several industries, but have actual different meanings like discussed below. This era probably peaked around the IBM Watson victory in jeopardy, but with way less interesting results in healthcare [182], and 2010s brought deep learning. Nowadays, AI is a buzzword that is used to describe a wide range of systems, from the most simple to the most complex. It is clearly trending, reports on AI show that papers regarding the subject have seen a 20-fold increase from 2010-2019. For defining AI we could use the definition provided by a group of experts the European Commission asked to

write some guidelines on AI [51] From this document, it is understood, that first, we need to address the difference between intelligence and rationality. Since the first is more subjective and even philosophical, the second is more pragmatic and related to the capability of choosing the best action to take in a certain scenario towards a certain goal. This is a more concrete concept and although it is not the same as intelligence, it should be a part of it [51, 170]. From these two concepts, we can go even deeper, and define that rationality can be achieved in a AI system by perceiving the environment, reasoning with what is perceived, and acting on the environment. From these three elements, we can argue that reasoning is the core functionalit, which is related to taking data, understanding it or interpreting it, and reasoning on this data through a model (numerical or symbolical) to reach the best action.

There is also the need to address the current distinction for AI which is the narrow and general AI. The first is the one that exists nowadays and its a AI that is not generic, it is focused on a specific task. The second is the one that is not yet achieved, and it is the one that is more generic and can be applied to several tasks. This is the one that is usually associated with the popular or common concept of AI [51, 170]. So, with this is mind, we reached the definition of AI as:

*Artificial intelligence (AI) systems are software (and possibly also hardware) systems designed by humans that, given a complex goal, act in the physical or digital dimension by perceiving their environment through data acquisition, interpreting the collected structured or unstructured data, reasoning on the knowledge, or processing the information, derived from this data and deciding the best action(s) to take to achieve the given goal. AI systems can either use symbolic rules or learn a numeric model, and they can also adapt their behaviour by analysing how the environment is affected by their previous actions. As a scientific discipline, AI includes several approaches and techniques, such as machine learning (of which deep learning and reinforcement learning are specific examples), machine reasoning (which includes planning, scheduling, knowledge representation and reasoning, search, and optimization), and robotics (which includes control, perception, sensors and actuators, as well as the integration of all other techniques into cyber-physical systems). [51]*

Of course, since the developments of this area have been so vast, this concept may become outdated very quickly. However, it is a good starting point to understand the concept of AI and its implications.

## 2.2 Evidence Based Medicine

In 1996, David Sacket and colleagues defined EBM as *the conscientious, explicit, and judicious use of current best evidence in making decisions about the care of individual patients* [171]. Despite having historical antecedents dating back to at least the 19th century, the first time the term "evidence-based medicine" was first coined by a team at McMaster University in Canada in the 1980s [184]. This was a time when clinical decision-making was mostly based on untested observations and physicians' experience, leading to variability in treatment strategies. The birth of

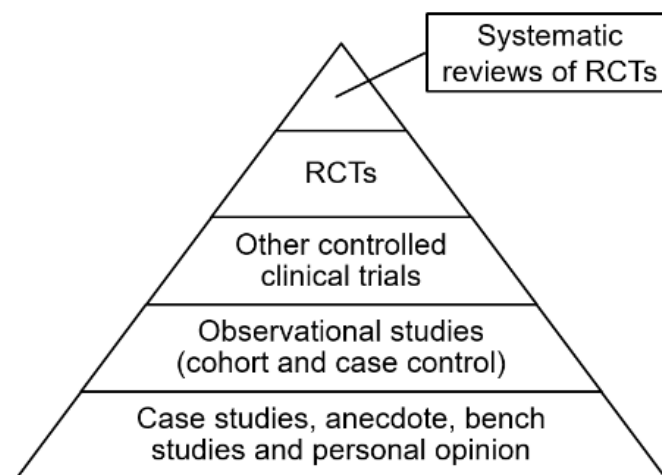


Figure 2.1: EBM adapted from [79]

EBM marked a pivotal moment in medical history, aiming to standardize patient care and improve outcomes. So, EBM is still a relatively recent concept in healthcare, which entails integrating the best available research evidence with clinical experience and patient values to make decisions about patient care. With this, we can, as stated by Sacket, define EBM into 3 major pillars:

- Best available evidence
- Clinical expertise
- Patient values, expectations, and/or wishes

Clinical expertise refers to the acumen and discernment gained from hands-on clinical experiences and consistent practice. This expertise manifests notably in enhanced diagnostic abilities and in the considerate recognition of a patient's unique circumstances, rights, and wishes when making care decisions. The term 'Best available evidence' pertains to pertinent clinical studies, often stemming from epidemiological investigations. This is linked with the ability (and willingness) to challenge current diagnostic methods and treatments, introducing alternatives that are more robust, precise, effective, and safer. Without experience, clinical practices blindly follow the best available evidence, which is not always the best option for the patient, since sometimes it may be inapplicable to a specific scenario. Without evidence, clinical practice becomes stagnated and unable to evolve [171].

The main concept of EBM is the hierarchy of evidence, which classifies different types of research studies based on their methodological quality and applicability to patients. At the top of this hierarchy are RCTs and systematic reviews of RCTs, which are considered to provide the most robust evidence. Observational studies, case series, and expert opinions are further down the hierarchy due to their inherent limitations (figure 2.1). EBM advocates for the application of the highest level of evidence available in clinical decision-making.

Historically, medical decisions leaned heavily on anecdotal observations and the prevailing beliefs of seasoned practitioners. To underscore the dangers of relying solely on such expert opinions, Sackett frequently recounted the circumstances surrounding George Washington's un-

fortunate end. Despite being in good health at the age of 68, Washington developed epiglottitis. Rather than opting for a tracheostomy, a treatment method known since ancient Greek times, his physicians, guided by the prevailing expert opinion, chose bloodletting as the course of action. Tragically, this decision led to Washington's likely preventable death, highlighting the critical importance of grounding medical decisions in robust evidence. However, EBM is not without critiques. The first one is that this is what medicine is all about and is already practiced all over. The data suggest something different [171]. The second refers to the virtually impossible task of keeping up with the literature. This argument, despite being refuted by examples of clinicians doing it, does raise the question of how can we deal with this, taking into account the increasing evidence overflow that the current times bring. How can we keep up with the literature and how can we make sure that the evidence is being applied in clinical practice? This is a very important question since the evidence is only useful if it is applied. This is where KDD and AI can play a role, as we will see in the next sections.

## 2.3 Extracting Knowledge of Data

KDD is about turning data into knowledge. However, turning data into knowledge or insights is not new in healthcare. The first attempts to use data to improve healthcare date back to the 17th century, when John Graunt used data from the London Bills of Mortality to study the causes of death in the city [3]. This was the first time that data was used to understand the health of a population. Since then, the field of KDD has evolved significantly, and it is now a crucial part of healthcare, helping to improve patient outcomes, enhance clinical decision-making, and optimize healthcare delivery. Additionally, the fact that data is being collected at an unprecedented rate, and the need to extract knowledge from it, has led to the development of several methodologies and frameworks to map low-level data (granular) into short reports, more abstract or more useful formats [69]. So it is only natural to see that KDD has become very popular in a wide range of industries nowadays. Healthcare is no exception and KDD has been applied to several areas of healthcare, from clinical decision support to disease surveillance and outbreak detection. Reports and papers suggest that [57] the digital data in the healthcare space has been increasing rapidly, due to the adoption of EHR and similar digital tools in the healthcare space. The complexity and vastness of healthcare data, encompassing electronic health records, genomic data, medical imaging data, and various other types of data, call for the adoption of intelligent systems that can mine this data for useful insights. The KDD process, comprising data cleaning, integration, selection, transformation, data mining, pattern evaluation, and knowledge presentation, can effectively help discover patterns and relationships in healthcare data, which are often not apparent to traditional analysis methods. This process facilitates the prediction of disease outbreaks, the identification of high-risk patient groups, the optimization of treatment plans, and the enhancement of healthcare service delivery. The generic process for KDD is shown in figure 2.2.

Several frameworks have been proposed to implement the KDD process. One such prominent framework is Cross-Industry Standard Process for Data Mining (CRISP-DM), which comprises



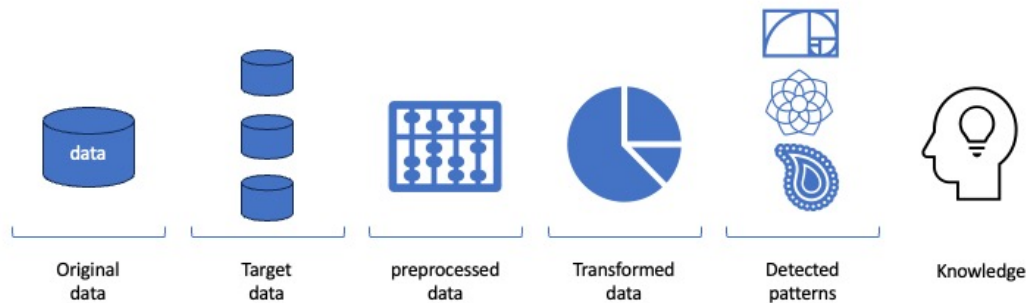


Figure 2.2: KDD Process, adapted from [69]

business understanding, data understanding, data preparation, modeling, evaluation, and deployment. CRISP-DM was conceived in 1996 and became a European Union (EU) project under the ESPRIT funding initiative in 1997 [38]. Sample, Explore, Modify, Model, and Assess (SEMMA) [165] involves five stages: sampling, exploration, modification, modeling, and assessment. It starts by analyzing a subset of data, then seeks patterns and modifies variables. A model is built and the results are evaluated. While SEMMA covers key data-mining aspects, it misses fundamental components of information system projects like analysis and implementation.

It is important to distinguish however that KDD is not the same as Data Mining. Like stated in [69], we agree that KDD is a major process of which Data Mining is a part. So, in order to understand the process of KDD, we need to understand the process of Data Mining, which can be understood as the application of algorithms for extracting patterns from data. There are several classes of algorithms, each best suited for different kinds of tasks:

- **Classification Algorithms:** These are used to predict categorical class labels. Examples include Decision Trees, Naive Bayes, Support Vector Machines (SVM), K-Nearest Neighbours (KNN), and various types of Neural Networks. These are used in disease diagnosis, patient risk prediction, and readmission prediction.
- **Clustering Algorithms:** These are unsupervised methods used to group similar data points together. K-Means, Hierarchical Clustering, DBSCAN, and Self-Organizing Maps are common clustering algorithms used in patient segmentation and anomaly detection.
- **Regression Algorithms:** These are used to predict continuous output variables. Examples include Linear Regression, Logistic Regression, and Regression Trees. These algorithms find application in predicting disease progression and healthcare costs.

- **Association Rule Mining Algorithms:** These discover associations or patterns among a set of items in large databases. Apriori and FP-Growth are commonly used algorithms in this class, helping in discovering co-occurring health conditions or drug interactions.
- **Sequential Pattern Mining Algorithms:** These help discover or predict specific sequences of events, which is particularly useful in medical trajectory analysis.
- More sophisticated architectures and algorithms appeared with neural networks, generative AI, and reinforcement learning, among others

As a result, KDD is the process of applying Data Mining algorithms to data but also the data preparation, selection, cleaning, and most important of all, the incorporation of prior knowledge about the domain along with the proper interpretation of results. This difference is vital to understanding KDD since blindly applying data mining or ML methods to data will only render results that are not useful or even misleading [69].

In short, KDD can be understood as a multidisciplinary subject that bridges and aggregates knowledge from different areas like ML, pattern recognition, databases, statistics, AI, knowledge acquisition for expert systems, data visualization, and high-performance computing. On top of all of these subject and research areas, sits the most important of all - which is domain expertise.

## 2.4 Health Data Science

Health Data Science is an interdisciplinary field that applies rigorous methods to transform health-care data into actionable knowledge for improving health outcomes. It involves the collection, interpretation, and application of vast amounts of biological, clinical, population, and health system data to improve patient care and public health. The advent of electronic health records, genomics, mobile health technologies, and other forms of big data have fueled the growth of this discipline.

In practice, Health Data Science involves the use of statistical and machine learning methods to analyse healthcare data. This data can be patient records, genomic data, demographic data, and more. It includes elements from various disciplines like biostatistics, epidemiology, informatics, and health economics. The ultimate goal is to provide a data-driven foundation for health decision-making for clinicians, health administrators, policymakers, and researchers.

An integral part of Health Data Science is predictive modeling and hypothesis testing. Predictive modeling involves the creation and use of statistical models or machine learning algorithms to predict future outcomes based on historical data. Hypothesis testing, on the other hand, is used to test the validity of a claim or theory about a population based on sample data. These are crucial for health data science as they allow us to make educated guesses about health trends and outcomes.

Importantly, Health Data Science has significant ethical and privacy considerations. Health data is often sensitive and personal, so maintaining privacy and confidentiality is crucial. This requires secure data handling and storage practices, as well as careful consideration of ethical implications when designing studies and algorithms. Health Data Scientists must also be wary of algorithmic bias and must ensure their models do not perpetuate or amplify health disparities. The

ultimate goal of Health Data Science is to improve patient outcomes and health equity using the best available data and methods.

The potential of using systematically created data in healthcare has certainly a lot of potential. However, we have seen in the past as well, that the hype of AI and ML usually are not supported by truth. There are currently six main aspects that hinder the potential of health data science [141, 147]:

- interoperability
- semantic
- secondary usage
- data quality
- privacy and ethics
- observational data

**Interoperability** is defined by *the ability of two or more systems or components to exchange information and to use the information that has been exchanged*[1]. In the context of healthcare, this means that different systems should be able to exchange data and interpret the data that has been exchanged. This is a very important aspect of health data science since the data is usually stored in different systems, with different structures and different purposes. So if systems are locked inside themselves and no export is possible, data becomes inaccessible. So, it is only natural that interoperability has been a key factor in gathering data. With tens or hundreds of different systems in every health institution, the possibility of exchanging data between EHRs plays a vital role. The usage of interoperable standards is of extreme importance in order to tackle the need to get data with a predefined structure.

**Semantic** adds a layer to the previous points, being sometimes related to interoperability as well. The fact that several institutions and EHRs are involved in creating knowledge from data, raises the problem that not all have data coded in clinical terminologies, or if they do, it is seldom the same across systems, since semantics has a very tight relationship with domain, especially in healthcare. So the normalization of uncoded terms is often required and mapping across terminologies is also very common, which is time-consuming and requires expertise in several fields.

**Secondary usage** is related to the fact that we are aiming to use data for a purpose for which the data was not created. The main goal of the healthcare data is to provide care. It is not meant for analysis and gaining insights. More than that, is already pretty well documented that the usage of EHR is very different from institution to institution and from country to country [12, 209, 147]. This means that the context where data was collected, even the actual person who inserted the information could be key to interpreting the results. To make things more complicated, the degree of precision of the data inserted varies highly on the type of information and context, as reported in [56].

**Data Quality** stems from the secondary usage. If the data is not reliable, how can we use it to gather useful knowledge from it? In order to, at least, try to counter this, we can apply several statistical methods and machine learning algorithms to try to clean the data. However, this is not a trivial task, since the data is usually very heterogeneous and the context where it was collected

is not always available. So, data quality is a very important aspect of health data science, since it can be the difference between a good and a bad model.

**Privacy and ethics** adds yet another layer to the problems of health data science. The fact that we are dealing with sensitive private data, which is not meant to be used for secondary purposes, raises the question of privacy and ethical concerns. Anonymization techniques and privacy-preserving methods are key to tackling this problem. However, they are not problem-free and are often complicated to assess. Moreover, the risks are very high, since the data is very sensitive and the consequences of a breach of privacy can be very serious, undermining public trust in clinicians, healthcare institutions and the healthcare system as a whole.

**Observational Data** relates to the fact that all health data science will be based on observational data. This means that the data is not collected in a controlled environment, which is the case for RCTs. Consequently, this data is subject to several biases, which are not always possible to control. The cornerstone of RCTs is simply not possible to apply here, preventing a proper comparison between groups. Even though there are techniques to tackle the unbalance in the measured variables, there is no way to control the unmeasured variables, which can be the cause of the observed effect. This is of particular importance and a major area of research at the moment, as we will see in the sections 2.5 and 2.6.

With this in mind, it is natural to assume that health data science and EBM are very synergic. If, on the one hand, we could argue that KDD can take EBM even further by using Data Mining and AI to produce synthetic evidence by analyzing, summarizing, or even combining evidence from several sources in order to feed medical practice with the best evidence available in a useful manner. On the other hand, we could also argue that EBM can be used to guide the KDD process, by providing the necessary domain knowledge to interpret the results and to guide the process of data preparation, selection, and contextualization. The domain knowledge mentioned in the KDD section could be applied by EBM.

The synergy of KDD and EBM has the potential to revolutionize healthcare delivery and improve patient outcomes. By leveraging the power of data analysis and advanced algorithms, health data scientists can identify novel biomarkers, develop predictive models, and personalize treatment plans based on individual patient characteristics. This not only enhances clinical decision-making but also enables precision medicine, where treatments can be tailored to the specific needs of each patient. Additionally, the use of health data science in evidence-based medicine allows for the continuous monitoring of treatment effectiveness and safety, facilitating the identification of best practices and the refinement of clinical guidelines over time.

## 2.5 Explainable Artificial Intelligence

AI has experienced unprecedented advancements in the last decade, leading to its integration in various domains, including medicine. It has been instrumental in transforming clinical decision-making, drug discovery, patient monitoring, and predicting disease trajectories. Despite these advancements, the "black box" nature of complex AI models poses interpretability challenges,

limiting their widespread adoption in healthcare, a field where transparency, reliability, and understanding of decision-making processes are vital. This lack of interpretability, also known as opacity, can lead to misdiagnoses, inappropriate treatment plans, and, most importantly, breaches in trust among clinicians, patients, and AI systems.

As such, the concept of Explainable AI (XAI), which aims to create a suite of techniques that produce more explainable models while maintaining a high level of predictive accuracy, has gained significant attention in medical AI research. XAI seeks to bridge the gap between AI opacity and human interpretability, and in doing so, it can enhance the transparency, reliability, and acceptance of AI applications in the healthcare setting.

So, for this to happen, we need a new framework for applying such mechanisms. A new step that could be attached to the ones seen before in section 2.3 will enable human comprehension of the model's output.

Even though several grouping and taxonomies of XAI are available mentioned in [7, 114, 23, 114, 101], a simplified approach based on [101] will be used in order to contextualize this concept.

We can divide it into two main categories. Firstly the explanation type is divided into global and local. Local and global explanations are methods used to interpret machine learning models, especially those that are considered "black box" models, such as deep learning networks. These methods help us understand why and how a model makes certain decisions, which can be crucial in many settings for ethical, legal, and practical reasons.

**Local Explanations:** These involve understanding the prediction of a ML model for a specific individual instance. They help to answer questions like: "Why did the model predict that this particular patient has cancer?" or "Why was this specific transaction flagged as fraudulent?".

**Global Explanations:** These focus on understanding the model behavior across all instances, or more broadly on a dataset-wide level. They help to answer questions like: "What features are generally important for prediction in the model?" or "What is the overall logic of the model?".

Secondly, we have the method type, where we have 3 main subcategories related to the stage of the data science process it is applied, *pre*, during, and *post*-model training.

**Pre-Model XAI:** These methods involve improving the transparency and interpretability of models before they are even trained. This includes thoughtful feature engineering, Exploratory Data Analysis (EDA), and applying domain knowledge to create meaningful variables. The goal is to design a model that will be more interpretable from the onset.

**Intrinsic XAI:** This involves using machine learning models that are intrinsically explainable. These models are designed in such a way that their decision-making process is understandable by default. Examples include linear and logistic regression, cox regressions, decision trees, Naïve Bayes, Bayesian Network (BN), and rule-based models. While these models may sometimes lack the predictive power of more complex models, they provide clear interpretability: you can directly examine the impact of the variables and understand how the model makes its predictions.

**Linear Regression** is a linear approach to modeling the relationship between a dependent variable and one or more independent variables. It assumes that the relationship between these variables is linear and can be represented by a straight line. The goal is to fit the best possible line that

describes this relationship by minimizing the sum of the squared differences (errors) between the observed values and the values predicted by the line. Linear regression is widely used in various fields for prediction, modeling, and determining the strength and character of the relationship between variables. It forms the basis of many more complex statistical modeling techniques.

**Logistic Regression** is used to model the probability of a binary outcome that depends on one or more independent variables. Unlike linear regression, which predicts a continuous outcome, logistic regression predicts the probability of a categorical outcome (e.g., success/failure, yes/no, 1/0). The logistic function is applied to the linear combination of independent variables to ensure that the estimated probabilities are between 0 and 1. It's often used in fields like medicine, economics, and social sciences to predict the likelihood of an event occurring based on various factors.

**Cox Regression** or the Cox proportional-hazards model, is a statistical technique used for investigating the effect of several variables on the time a specified event takes to happen. In medical research, this often refers to survival times. The model allows for the estimation of hazard ratios, which describe how the hazard changes with a one-unit change in the predictor variable. The Cox model makes an assumption that the hazard ratios are constant over time, known as the proportional hazards assumption. This model is vital for understanding how different factors influence survival or failure time and is commonly applied in epidemiological and medical research.

**Bayesian Networks** A BN, also known as a belief network or Directed Acyclic Graph (DAG) model, is a probabilistic graphical model that represents a set of random variables and their conditional dependencies via a DAG.

Given a set of variables  $X = \{X_1, X_2, \dots, X_n\}$ , the joint probability distribution is given by:

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Parents}(X_i))$$

where  $\text{Parents}(X_i)$  is the set of parent variables of  $X_i$  in the network.

This formula represents the factorization of the joint distribution over  $X$ , based on the graphical structure of the Bayesian network.

Now, in the Bayesian network, each node is conditional independent of its non-descendants given its parents. If we denote  $ND(X_i)$  as the set of non-descendants of  $X_i$  and  $Pa(X_i)$  as the parents of  $X_i$ , the conditional independence is described as:

$$X_i \perp ND(X_i) | Pa(X_i)$$

This means that  $X_i$  is conditionally independent of its non-descendants given its parents.

A common task for Bayesian networks is inference, which means computing the posterior probability of a set of query variables  $Q$ , given some observed variables  $E$ . That is, we want to compute  $P(Q|E)$ . According to the Bayes rule, we have:

$$P(Q|E) = \frac{P(Q, E)}{P(E)} = \frac{P(Q, E)}{\sum_{q \in Q} P(Q = q, E)}$$

where the denominator is a normalization constant ensuring the result is a valid probability distribution. Note that performing this inference is NP-hard, which is why various approximation algorithms have been developed.

**Tree based methods** Tree-based machine learning methods are a subset of algorithms that use a tree-like graph structure for making decisions or predictions. The most basic type is the Decision Tree, where the tree is used to go from observations about an item to conclusions about the item's target value (classification or regression). Each node in the tree represents a feature in the dataset, each branch represents a decision rule, and each leaf node represents the output value. More advanced tree-based methods include Random Forests, which build multiple Decision Trees and average their predictions for better accuracy and generalization, and gradient-boosted trees, which build trees sequentially, each one correcting the errors from the previous one.

The major advantage of tree-based methods is their ease of interpretation and understanding, especially for Decision Trees. However, a single tree is often prone to overfitting, where it performs well on the training data but poorly on unseen data. This is why ensemble methods like Random Forest and Gradient Boosting are popular; they aim to increase robustness and predictive power by combining multiple trees. These methods are widely used in various domains including but not limited to finance, healthcare, and natural language processing for tasks like classification, regression, and even unsupervised learning tasks like clustering.

**Post-Hoc XAI:** Post-hoc methods are applied after a model has been trained, to try to explain its decisions. This includes techniques like feature importance analysis, partial dependence plots, Local Interpretable Model-Agnostic Explanation (LIME), SHapley Additive exPlanations (SHAP), and counterfactuals. For instance, LIME can be used to create local explanations for individual predictions made by any model, and SHAP values can be used to interpret the impact of features on the model's output both locally and globally. Counterfactuals try to explain a model by example, providing possible changes that would alter the outcome provided by the model.

It is to be noted that a methodology can be classified into two categories. For example, LIME is a local explanation model in a *post-hoc* manner.

Despite all of this, we have to take into account that pre-model and post-hoc methodologies are a proxy for an explanation of the models. That is why we could argue, as stated in [168] that only an intrinsically transparent model can really be the basis of XAI. While *post-hoc* or pre-model methods are only a potentially unreliable proxy for an explanation.

## 2.6 Causality

Using once again the tale of George Washington, but now with a different purpose, the medical doctors in that region of the globe at least, followed the theory of humors which relied on the fact the healthy human person was a balance between 4 humors (blood, phlegm, yellow bile, and black bile). So, the treatment for George Washington was to rebalance those 4 humors, and so, doctors needed to remove blood, which was the supposed cause of his illness. Since microbiology and its importance would only be discovered sometime after, the idea at the time was inspired by



the fact that the imbalance of these 4 senses of humor and illness present at the same time. This is now known, as a textbook definition of confounding correlation with causation. And in this subject in particular, it was not the imbalance in the humors that caused illness, but an illness that caused the imbalance. So, this example shows that evidence without proper causality can lead to misguided results and mistrust. So that is why, nowadays, EBM, XAI can be brought together and expanded through Causal Machine Learning (CausalML). But what is causality? We could argue that is related to something **causing** something else. This causal effect, especially in medicine can be related to a comparison of the outcome a particular person would exhibit given a particular intervention and the outcome in the same person of the control intervention. This is particularly hard since we cannot do both things at the same time. This is the base of why RCTs are the gold-standard of experimentation, since they are the current best tool to achieve something similar to this [166].

CausalML is a branch of machine learning that focuses on understanding and quantifying causal relationships from data [87]. Instead of just finding patterns or correlations in data, CausalML aims to uncover the cause-and-effect relationships that explain these patterns. This is especially important since current or traditional ML and AI methodologies rely heavily on association and not causation. So, CausalML can support traditional algorithms to solve its limitations [145]. There are currently two main frameworks for trying to unveil causality in data: the Structural Causal Model (SCM) and the Potential Outcome Framework (POF) [175]. **SCM** relies on 1) Causal Graphs and 2) structural equations.

1. Causal Graphs are based on DAGs: These are graphical models used to represent causal relationships between different variables. The nodes in the graph represent variables, and the edges (arrows) between nodes represent causal relationships. For instance, an edge from Node A to Node B signifies that A has a causal effect on B. We should not confuse causal graphs with BN. Even though both rely on DAGs, Causal Graphs represent causal relationships, and BN represent conditional dependencies.
2. Structural Equations refer to a set of mathematical expressions that represent causal relationships between variables. These equations model the way changes in one variable, often termed the "cause," lead to changes in another, termed the "effect." Within a structural equation model (SEM), both observed and latent (unobserved) variables can be incorporated, and the causal pathways between them are explicitly defined. By employing SEM in CausalML, researchers can elucidate intricate relationships among variables, disentangle direct from indirect effects, and infer causal mechanisms. This approach provides a more profound understanding of the underlying data-generating process, enabling better predictions and interventions in complex systems.

**POF** model centers on the concept of potential outcomes which can be understood as all of the possible outcomes for a patient. Each unit (e.g., a patient or a sample) has a set of potential outcomes, each corresponding to one of the possible treatments the unit could receive. The causal



effect is defined as the difference between these potential outcomes. This framework allows for the formal definition and estimation of causal effects. In this approach, we consider the potential outcomes for each unit (for example, a patient in a healthcare context) under each possible treatment or intervention. Each unit has a set of potential outcomes corresponding to each possible intervention. However, we can only observe one of these outcomes for each unit, corresponding to the intervention that was actually received. The other outcomes, which would have occurred had different interventions been implemented, remain latent. These are known as counterfactual outcomes.

The difference between potential outcomes under different treatments represents the causal effect of the treatments. For instance, in a healthcare scenario, if we are studying the effect of a drug, we might consider two potential outcomes for each patient: the outcome if the patient is given the drug, and the outcome if the patient is not given the drug. The difference between these outcomes represents the causal effect of the drug on the patient. However, as we can only observe one of these outcomes for each patient (the one corresponding to the treatment they actually received), a key challenge in causal inference is estimating the unobserved potential outcomes. Various statistical methods, including randomized experiments, matching methods, and instrumental variable methods, can be used to estimate these unobserved potential outcomes.

1. **Counterfactuals:** This is a concept rooted in the idea of "what-if" scenarios. A counterfactual outcome for a given individual is the outcome that would have occurred had the individual been exposed to a different treatment or condition. Counterfactuals play a pivotal role in the field of causal machine learning, offering a sophisticated approach to understanding cause-and-effect relationships. In essence, a counterfactual is a conceptual device used to contemplate what would have happened under a different set of circumstances than what actually occurred. This hypothetical scenario is created by altering some aspect of the actual situation, providing a means of comparison to evaluate the effect of a particular variable or intervention.

For instance, in the context of healthcare, consider a scenario where a patient was given a particular drug and recovered. The counterfactual question here would be: "What would have happened to the patient if they hadn't been given the drug?" Answering this question allows us to estimate the causal effect of the drug on the patient's recovery. While the true counterfactual outcome is unobservable (since we cannot rewind time and alter the decision), various statistical techniques, machine learning algorithms, and experimental designs are employed in causal inference to estimate this effect as accurately as possible. The ability to make such counterfactual inferences is crucial in numerous fields, including medicine, economics, social sciences, and policy-making, where understanding causal relationships is paramount.

2. **Instrumental Variables:** These are variables that are related to the treatment but not the outcome, except through their effect on the treatment [33, 59]. They can be used to control for unmeasured confounding variables. Instrumental variables (IVs) are a powerful tool

used in causal inference to help address the problem of confounding variables, especially in situations where randomization is not feasible. An instrumental variable is a variable that is correlated with the independent variable (the treatment) but does not directly affect the dependent variable (the outcome), except through its effect on the treatment. In other words, it is a variable that induces changes in the explanatory variable but is otherwise unrelated to the outcome of interest.

The idea behind using an instrumental variable is to isolate the portion of the variation in the treatment that is independent of the confounders and therefore provides a "natural" form of randomization. The causal effect of the treatment on the outcome can then be estimated based on this variation.

For example, in a study assessing the impact of education on income, it's challenging to identify causal effects because numerous unobserved factors (like ability or motivation) could affect both education and income, thus confounding the relationship. If we find an instrumental variable – say, distance to the nearest college (which affects the likelihood of getting higher education but doesn't directly affect income) – we can use this to isolate the part of the variation in education that is unrelated to the unobserved confounders, and thereby get a more accurate estimate of the causal effect of education on income.

It's crucial, however, to remember that the use of instrumental variables relies on certain assumptions, such as the relevance and exogeneity of the IV. The relevance assumption requires that the IV is correlated with the treatment, and the exogeneity assumption requires that the IV affects the outcome only through the treatment and is not related to the unobserved confounders. Violations of these assumptions can lead to biased and inconsistent estimates of causal effects.

3. **Propensity Score:** This is the probability of a unit (e.g., a patient) being assigned to a particular treatment given a set of observed characteristics. Propensity scores are used to balance the characteristics of treatment and control groups, mimicking the conditions of a randomized experiment [14, 16]

The propensity score is a statistical concept widely used in causal inference, particularly in the field of observational studies where random assignment of treatment is not possible. The propensity score for an individual is the probability of receiving the treatment given the observed characteristics of that individual. In other words, it's the likelihood that a particular individual would be assigned to the treatment group based on their observed features.

The key idea behind propensity scores is to create a balance between the treatment and control groups based on these observed characteristics, thus mimicking the conditions of a randomized controlled trial. This balance helps to eliminate bias caused by confounding variables, allowing for a more accurate estimate of the treatment effect. Once propensity scores are calculated, they can be used in several ways including matching, stratification,

Inverse Probability of Treatment Weighting (IPTW), and as covariates in regression adjustment.

For example, consider a study investigating the effect of a training program on job outcomes. Individuals might self-select into the training program based on characteristics like motivation or prior education, which are also related to job outcomes, creating confounding. The propensity score, calculated based on these observed characteristics, can be used to match each participant in the training program with a similar non-participant or to weight the observations, such that the distribution of observed characteristics is similar between the groups. This helps to isolate the effect of the training program on job outcomes.

After achieving this balance, it becomes more meaningful and less biased to estimate treatment effects, such as Average Treatment Effect (ATE) and Average Treatment Effect on the Treated (ATT).

The ATE quantifies the difference in mean outcomes between units that are treated and units that are not. Essentially, it calculates the expected difference in outcomes if everyone in a population received a treatment versus if no one received it. Mathematically, the ATE is represented as:

$$ATE = E[Y_1 - Y_0]$$

where  $Y_1$  is the potential outcome under treatment and  $Y_0$  is the potential outcome under control. The expectation is taken over the entire population.

After addressing confounding using propensity scores, the ATT narrows its focus to the treated subpopulation. It measures the average effect of a treatment on those units that actually received the treatment, comparing their observed outcomes to what their outcomes would have been without the treatment. The formula for ATT is:

$$ATT = E[Y_1 - Y_0 | D = 1]$$

where  $Y_1$  and  $Y_0$  once more denote potential outcomes under treatment and control, respectively, and  $D$  is an indicator for treatment (with  $D = 1$  indicating treatment).

However, it's important to note that propensity scores only account for observed confounders. If there are unobserved confounders that influence both treatment assignment and the outcome, propensity score methods may still produce biased estimates of the causal effect.

## 2.7 Legal and ethical considerations

As Health Data Science, KDD and AI in healthcare get more and more popular, it is important to consider the words postulated by Francis Bacon in the "Wisdom of the Ancients", "*mechanical arts are of ambiguous use, serving as well for hurt as for remedy.*" [2]. This is currently as true for AI as it was at the time. We must consider the good and the bad of such technologies, and how to mitigate the bad and enhance the good. In this section, we will discuss the legal and ethical

considerations of AI in healthcare. Ensuring the proper use of healthcare data is key to preserving public trust and ensuring the long-term viability of data-driven health initiatives.

One of the primary legal considerations is data privacy. Laws such as the Health Insurance Portability and Accountability Act (HIPAA) in the USA, and the General Data Protection Regulation (GDPR) in the EU, set stringent rules on how healthcare data should be stored, shared, and processed. They require data scientists and healthcare providers to take steps to anonymize data and limit the scope of data usage. Breaching these regulations can lead to severe penalties, including fines and imprisonment. Secondly, there's the matter of data security. With the rise of cyber-attacks, ensuring the robustness of the system against such breaches is both a legal requirement and an ethical obligation. Security breaches could lead to sensitive patient data being stolen, with severe implications for the individuals involved and for the trust in the healthcare system as a whole.

The European Health Data Space (EHDS) refers to a strategic initiative by the EU aimed at creating a unified and secure platform for sharing and accessing health-related data across member states. AI is expected to have a significant impact on the EHDS in several ways [68]:

- **Improved Diagnostics and Personalized Medicine:** AI can analyze vast amounts of health data, including medical records, imaging, and genetic information, to enhance diagnostic accuracy and tailor treatments to individual patients. This can lead to more effective and efficient healthcare delivery.
- **Data Integration and Interoperability:** AI can help harmonize data from various sources within the EHDS, including electronic health records, wearable devices, and clinical databases. This promotes interoperability, allowing healthcare professionals to access comprehensive patient information seamlessly.
- **Predictive Analytics:** AI-powered predictive models can help forecast disease outbreaks, patient admission rates, and healthcare resource utilization. This enables better resource allocation and proactive healthcare planning.
- **Drug Discovery and Development:** AI can accelerate drug discovery by analyzing genetic data, identifying potential drug candidates, and predicting their efficacy and safety profiles. This can expedite the development of new treatments and therapies.
- **Enhanced Clinical Decision Support:** AI can provide healthcare providers with real-time decision support, offering recommendations based on the latest medical evidence and patient-specific data. This can lead to more informed clinical decisions and better patient outcomes.
- **Data Security and Privacy:** The EHDS must ensure the privacy and security of health data. AI can help by implementing robust encryption, access controls, and anomaly detection systems to safeguard sensitive information.

- **Research and Insights:** AI can facilitate large-scale data analysis for medical research, enabling researchers to identify patterns, correlations, and potential breakthroughs in healthcare. This can lead to advancements in medical knowledge and treatments.
- **Patient Engagement and Monitoring:** AI-driven apps and wearable devices can empower patients to take a more active role in managing their health. These technologies can monitor vital signs, offer health advice, and send alerts to healthcare providers when necessary.
- **Reduced Healthcare Costs:** By optimizing healthcare processes, improving diagnosis accuracy, and preventing medical errors, AI can contribute to cost savings within the healthcare system, making it more sustainable.
- **Regulatory Challenges:** Implementing AI in healthcare requires navigating complex regulatory frameworks, ensuring ethical use, and addressing issues related to bias and fairness in AI algorithms. The EHDS will need to establish guidelines and standards to address these challenges.

On the ethical front, considerations include ensuring data fairness and avoiding bias. Given the diversity of patients in terms of age, race, sex, socioeconomic status, etc., algorithms should be designed and validated to ensure that they don't unintentionally perpetuate or amplify societal biases. For instance, a predictive model for disease risk should not unfairly disadvantage certain demographic groups. If we use data to derive knowledge and create Clinical Decision Support Systems (CDSSs) that orient and support clinical practice, they can be biased by the type of data that originated said knowledge [52, 22].

The importance of ethics in AI cannot be overstated, primarily because the decisions that these systems make can have profound implications on individuals and society. These decisions may affect anything from employment opportunities to legal outcomes, and increasingly, health outcomes. As AI models grow in complexity and application, they possess an enormous power that needs to be harnessed responsibly. This necessitates rigorous ethical considerations to ensure fair, unbiased, and transparent operations. Ethical lapses can result in discrimination, loss of privacy, and unjust outcomes, among other issues, which erode public trust in these technologies.

Equally important in the realm of AI is the understanding of why a model works the way it does. This concept, known as "explainability" or "interpretability", is central to AI ethics. It concerns the transparency of AI algorithms and the ability to understand and interpret their inner workings and decisions. Without this understanding, we run the risk of blind reliance on AI's 'black box' that may lead to erroneous or biased outcomes. It is critical to scrutinize AI models' reasoning processes, ensuring they align with human values and principles and are not based on inappropriate or discriminatory features.

In the context of healthcare, these considerations take on an even greater significance. AI applications in healthcare, such as diagnostic tools or treatment recommendation systems, directly impact human lives. They may influence critical decisions such as who gets treatment, what kind of treatment is administered, and when it should be given. These systems must not only be accurate

but also transparent, fair, and accountable. They should be designed and implemented in a way that respects patient rights, including privacy, autonomy, and informed consent.

Therefore, in healthcare, the need for ethical AI and model explainability is not just a matter of good practice, it's a matter of life and death. Bias or errors in AI could lead to misdiagnoses or inappropriate treatment recommendations, with potentially fatal consequences. Similarly, if AI-based systems make decisions that healthcare professionals or patients can't understand, it may lead to mistrust and potential harm. The advancement of AI in healthcare must ensure ethical considerations and explainability are at the core of AI model design, development, and deployment. This will build trust in AI systems and ultimately lead to better health outcomes.

Furthermore, the informed consent of patients is another significant ethical consideration. Patients should be fully informed about how their data will be used, and they should have the right to *opt-out* if they wish. Transparency is another crucial aspect that straddles both legal and ethical dimensions. It involves explaining how decisions or predictions are made by complex algorithms, particularly when they have significant implications for patient care. For instance, if an AI model is used to prioritize patients for treatment, it should be transparent about how the model makes its decisions. The explainability of machine learning models can help achieve this transparency, which aids in maintaining accountability and trust.

Finally, at the moment of this writing, there are in the EU several proposals that could impact AI in general and in healthcare in specific. The Medical Device regulation could impact the deployment of AI based systems and the AI act could also impact the development of AI based systems in healthcare.

*Nothing great in the world was accomplished without passion.*

Friedrich Hegel

# 3

## Case Studies

This chapter will comprise the work done during this PhD. The works developed and corresponding papers were a search for improving data usage in several steps of the KDD process. We can see some work dedicated to leveraging data acquisition or alternatives to it, like the works depicted in section 3.5, 3.6, 3.3, 3.2 and 3.1. Others will focus more on how to use the data in order to make a difference in clinical practice like the section 3.7, 3.8 and 3.4.

### 3.1 Can GANs help create realistic datasets?

This section is based on the paper entitled "GANs for Tabular Healthcare Data Generation: A Review on Utility and Privacy". It focuses on a review of the Generative Adversarial Network (GAN) framework for creating synthetic data for healthcare. Tries to compile the metrics used for comparing and assessing synthetic data in terms of utility - or how similar they are to the original data and privacy - how protective of the patient's data it is.

#### 3.1.1 Introduction

With the growing technological advances, the quantity of healthcare-related data produced around the world increased exponentially [45, 86]. Consequently, the potential for harvesting this data also increases. The value locked within this data could help provide better healthcare with new information about diseases, drugs, and preventive therapies. It can also help create better HISs, meaning an overall better clinical practice [21]. But for this to happen, data must reach capable hands at the right time. But the release of clinical data has several barriers attached and rightly so.

The leakage of patient's privacy can break the confidence of the population in healthcare professionals and institutions. Patient safety and privacy should be kept at all costs. However, the current mechanisms for privacy maintenance are very long, bureaucratic, and time-consuming, nationally [50], and internationally [137]. The current scenario and general methods for privacy safeguards are related to pseudo-anonymisation techniques. The removal of certain attributes, identifier modification, code grouping, or discretization are some methodologies. But not even these are totally safe [66]. Synthetic data appear as an alternative for clinical data sharing, promising great data utility with minimal privacy concerns. Synthetic data is data that is generated automatically through programmatic processes. This is especially impactful for the case at hand since synthetic data has no explicit connection with the original data. There are several mechanisms for data synthesis postulated by [76], there are process-driven methods and data-driven methods. Process-driven methods generate data through pre-determined models inputted into the generator. Data-driven methods produce new data based on inputted source data. With this, it is possible to create new patient data that has no relation to reality while providing the same statistical relations between variables. This provides the basis for quality clinical research on top of this new data. Even though these techniques are still new and in rapid development, the results seem interesting [76], but not without questions and doubts [178]. Creating a thorough survey based on the generation of synthetic data is seldom a simple task when compared to other surveys since synthetic data is present across several domains and has several uses, like software testing, assessing methods, or generating hypotheses. Moreover, synthesis has the double meaning of summing up information and generating something, easily wielding hundreds of results per query. Finally, trying to filter algorithms aimed at tabular data is also burdensome, since not always it is easy to discriminate input types. These factors make the survey interesting to focus on the state-of-the-art mechanisms of generating tabular data.

### 3.1.2 Theoretical background

First introduced in 2014, GANs [78] have been under the scope and have been proven very good for generating complex data. Images, text, and video have been successfully generated with very good performances. The original architecture is based on two artificial neural networks trained simultaneously in a competitive manner. One of them, the generator, has the objective of generating the most realistic possible data, while the second network – the discriminator, has the opposite aim of aiming to distinguish the realistic data from the synthetic data the best it can. So, the elegance of this architecture is that each network tries to make the other perform better every time. The GAN architecture is shown in 3.1.

The generator is represented by  $G_\theta$  where the parameter  $\theta$  represents the weights of the neural network. It takes as input, a Gaussian random variable, and outputs  $G_\theta(Z)$ . Distribution of  $G_\theta(Z)$  is denoted by  $P_\theta$ . The goal of the generator is to choose  $\theta$  such that the output  $G_\theta(Z)$  has a distribution close to the real data. The discriminator is represented by  $D_\omega$ , parametrized by weights  $\omega$ . The goal of the discriminator is to assign 1 to the samples from the real distribution  $P_X$  and 0 to the generated samples ( $P_\theta$ ). So, GANs can be mathematically represented by a minimax game



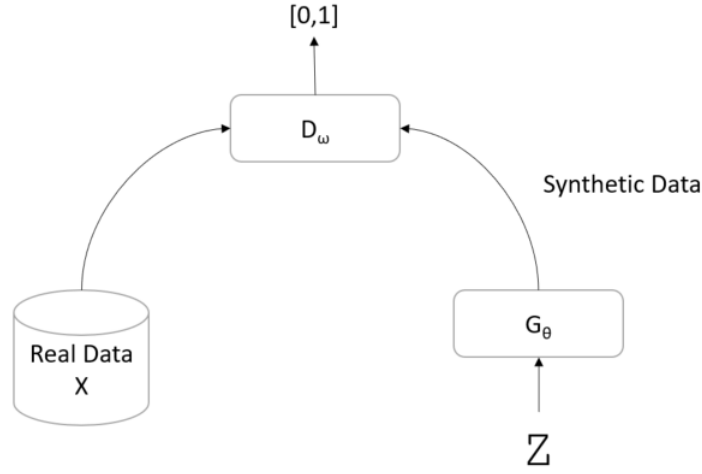


Figure 3.1: GAN framework

identified by:

$$\min_G \max_D E[\log(D_\omega(X)) + \log(1 - D_\omega(G_\theta(Z)))] \quad (3.1)$$

So,  $G$  must minimize this equation and  $D$  must maximize it, each one tweaking the weights of its network ( $\theta$  and  $\omega$ ) to do so. This is the loss function on the initial GAN architecture. After the classification of  $D$ , the  $G$  is trained again with the error signal from  $D$  through backpropagation. This equation is the log of the probability of  $D$  predicting that the real data is genuine and the log probability of  $D$  classifying synthetic data as not genuine. The equation is essentially the same as minimising the *Jensen-Shannon divergence* [78]:

$$\min_G JS(P_x || P_\theta) \quad (3.2)$$

Where the JS means the *Jensen-Shannon divergence* between the probability of the real data and the probability of the generated data. The JS divergence provides a measure of the distance between two probability distributions. Therefore, the minimization over  $\theta$  means, choosing the  $P_\theta$  that is closest to the target distribution  $P_X$  in the JS divergence distance. Despite the significant results provided by GANs with continuous real values, categorical values still seem to be a problem for this approach [109], since it is not directly applicable for calculating the gradients of latent categorical variables in order to train these networks through backpropagation. This happens since the output of the generator, even though can be transformed into a multinomial distribution with a softmax layer, sampling from it is not a differentiable operation, limiting the backpropagation process of the GAN.

### 3.1.3 Methods

This search was made between December 2020 and January 2021. It was made on “Web of Science”, IEEE, PubMed, Arxiv and finally GitHub. The terms searched were related to GANs, synthetic data generation, electronic health records, patient data, or tabular data. Applications of GANs to non-tabular data were filtered, like image, sound, video, or graphs. Time series and text data were also removed since the methodology for synthesizing this type of data has specific functions related to the nature of the data. The filter for date was after 2014 since GANs were introduced at that time. The queries used were similar to the one below, adapted for the search mechanics for each website.

("generation" OR "creation" OR "synthesis" OR "synthesizing" OR "generating" OR "creating") AND ("synthetic data" OR "synthetic patient" OR "synthetic electronic health record" OR "synthetic EHR" OR "realistic patient data" OR "realistic health record" OR ("synthetic" AND "privacy" AND "utility")) AND ("GAN" OR "Generative Adversarial Network")

From the total articles found (1165) with all the queries, 100 articles were chosen for full text and in the end, 22 papers with GAN implementations that were tested on tabular data were selected.

### 3.1.4 Results

The selected papers ranged from 2017 to 2020. Being that 2 are from 2017, 4 from 2018, 8 from 2019 and 8 from 2020. All authors showed original GAN implementations, apart from 2 papers. Beaulieu-Jones et al. [26] used a GAN architecture that was originally published with usage on image datasets [136]. Additionally, Vega-Marquez et al. [194] used an already known implementation of conditional GANs [126]. We classified papers regarding 3 metrics: utility, privacy and clinical. For utility, we looked for methods for measuring the generated data’s quality. As for privacy, we aimed for some mechanism for measuring the privacy loss of the new data. Concerning clinical metrics, any kind of evaluation from healthcare professionals was considered. This can be seen in table 3.1.

The metrics the authors used are exhibited in table 3.2.

Table 3.1: Summary of the articles selected.

ID	year	Acronym	Article	Metric	Code
1	2017	medGAN	[45]	Utility, Privacy, Clinical	[44]
2	2017	POSTER	[120]	Utility, Privacy	[150]
3	2018	table-GAN	[142]	Utility, Privacy	[128]
4	2018	dp-GAN	[214]	Utility, Privacy	[110]
5	2018	mc-medGAN	[35]	Utility	[34]
6	2018	TGAN	[219]	Utility	[154]
7	2019	PATE-GAN	[98]	Utility, Privacy	–
8	2019	SPRINT-GAN	[26]	Utility, Privacy, Clinical	[25]
9	2019	GAN-based	[119]	Utility, Privacy	–
10	2019	CTGAN	[217]	Utility	[153]
11	2019	WGAN-DP	[32]	Utility, Privacy	[31]
12	2019	PPGAN	[117]	Utility, Privacy	[116]
13	2019	medBGAN	[21]	Utility	–
14	2019	medWGAN	[19]	Utility	[18]
15	2020	ADS-GAN	[222]	Utility, Privacy	–
16	2020	corGAN	[189]	Utility, Privacy	[187]
17	2020	CGAN	[194]	Utility	–
18	2020	DPAutoGAN	[183]	Utility, Privacy	[64]
19	2020	GAN Boosting	[133]	Utility, Privacy	[132]
20	2020	RDP-CGAN	[190]	Utility, Privacy	[188]
21	2020	WCGAN-GP	[201]	Utility, Privacy	–
22	2020	SMOOTH-GAN	[160]	Utility	[65]

Table 3.2: Metrics utilised for evaluation

Acronym	Utility	Privacy
medGAN	1. Bern. 2. Pred F1	1. Attrib. disc. 2. Memb. inf. 3. KNN
POSTER	1. Pred Acc. 2. Corre. Mat. 3. BD	DP
table-GAN	1. Cumul. Dist. 2. Pred F1 MRE	1. Eucl. 2. Member. inf.
dp-GAN	1. Pred AUC 2. Bern.	DP
mc-medGAN	1. Pred F1 AUC 2. Bern. 3. ME F1 Acc	–
TGAN	1. KNN 2. NMI 3. Pred F1	–
PATE-GAN	1. Pred AUC AUPRC	DP
SPRINT-GAN	1. Pred AUC 2. Corre. Mat.	DP
GAN-based	1. Pred Acc. 2. Corre. Mat.	1. Hit. Rate 2. R. Linkage 3. Eucl.
CTGAN	1. Pred F1 R2 Acc.	–
WGAN-DP	1. Corre. Mat. 2. PCA 3. Pearson RMSE 4. Pred F1 RMSE 1-MAPE(F1)	1. Eucl. 2. Dupl. 3. DP
PPGAN	1. GS	DP
medBGAN	1. Assoc. Rul. 2. CCS Pred F1 3. KS	–
medWGAN	1. Assoc. Rul. 2. CCS Pred F1 3. KS	–
ADS-GAN	1. $\chi^2$ 2. JSD 3. WD 4. t-test 5. Pred AUROC 6. Corre. Mat.	DP
CorGAN	1. Pred F1 2. Bern.	Member. Inf.
CGAN	1. Pearson 2. Spearman 3. Pred F1 AUC Acc	–
DPAutoGAN	1. Pred AUROC R2 2. Bern.	DP
GAN Boosting	1. pRMSE 2. Pred AUROC AUPRC Acc.	DP
RDP-CGAN	1. Pred F1 AUROC AUPRC 2. MMD	DP
WCGAN-GP	1. Corre. Mat. 2. Pred F1	1. Dupl. 2. Eucl.
SMOOTH-GAN	1. DW MAE 2. Pearson 3. Pred AUROC AUPRC	–

Regarding privacy, 15 papers assessed it or included some kind of mechanism to improve data protection. The most common was including Differential Privacy (DP) in the generation process. Other mechanisms for measuring privacy loss were Membership Inference (Member. Inf.), Attributes Disclosure (Attrib. Disc.), Euclidean distance (Eucl.), record-linkage (R. Linkage) and Nearest Neighbours (KNN). As for utility, all papers assessed it. There were 3 major areas of utility assessment: Dimension-wise (DW) probability, cross-testing, and distance metrics. The most basic one was dimension-wise probability, which is important for making sanity checks for the generated data, comparing the distributions of each column between real and synthetic. In this category, we can find Bernoulli (Bern.), cumulative distributions (Cumul. Dist.), Pearson correlation (Pearson) and Spearman correlation (Spearman), correlation coefficients (CCS), chi-squared test ( $\chi^2$ ), *Kolmogorov-Smirnov* (KS) or Correlation Matrices (Corre. Mat.). Cross-testing was about training machine-learning algorithms with both datasets in order to compare the results. The key factor is generating a synthetic dataset based on the training set and then training models on the original training set and the generated dataset. Then the models are compared regarding their predictive capability on the (real) test set. This was a way of assessing if the generator models were capturing inter-variable relationships. The authors applied different metrics from AUC, F1, Area Under the Precision Recall Curve (AUPRC), Accuracy (Acc.) to Mean Relative Error (MRE). Finally, there was also the application of distance metrics, for measuring the difference between column distribution in both datasets. Jensen Shannon divergence (JSD), Wasserstein Distance (WD), Bhattacharyya Distance (BD) or Generate Scores (GS) that was a metric implemented by the authors of [117] that creates a metric based on the sum of the mean of *kullblack-leibler* distance of all columns. Other less used methods were Principal Component Analysis (PCA), Also, propensity score mean squared error ratio (pMSE). NMI (Normalised Mutual Information), which is the ability to capture correlations between columns by computing the pairwise mutual information and MMD (Maximum Mean Discrepancy), which is similar to distance metrics were also used. Regarding datasets utilized, the most used was MIMIC-III [97] (9 times). The papers used 27 different datasets, being 16 healthcare-related and 11 non-healthcare related. Finally, regarding clinical evaluation, only two papers assessed it, as it is possible to see in table 3.1. Both had a group of clinicians assessing a sample of both real and synthetic information and evaluating from 0 to 10, where 10 is the most realistic. One major point preventing a larger comparison is that despite some papers using the same dataset and same methodologies, the presented values are different, making it difficult for a clear comparison of results. One example is a dimension-wise prediction with F1 score for MIMIC-III. CorGAN presents the mean difference between the two classifications (real on real and synthetic on synthetic), while medBGAN presents the correlation coefficients of the two, and medGAN only presents the visual comparisons. Regarding code availability, 16 papers had the code publicly available in some form. As of January 2021, papers pointed in table 3.1 have public code.

### 3.1.5 Implications for future research

From the work done in this paper, it is clear that synthetic data generation is a growing field. The increasing number of papers through the years as the growing quality in the mechanisms of generating data and assessing its quality is clear proof. It also became apparent that privacy and utility in synthetic data represent a delicate balance. The very same definition of differential privacy represents it. The compromise between privacy and utility is real and should be taken into account when creating privacy-demanding datasets. Creating statistically good tabular datasets is already possible, but that task becomes increasingly difficult if privacy concerns are added. However, privacy is also a complex subject, and the context of the setting is important for privacy assessment, which explains the different approaches for evaluating privacy protection of synthetic data. From this review, we believe that a proper evaluation of synthetic data generators in the healthcare setting with privacy concerns should at least include utility and privacy evaluations. For utility, we believe that evaluating column-wise is a nice first check but insufficient alone. For table-wise, since there is no fundamental metric for assessing the inter-column correlations between mixed-type variables, cross-testing is the best next thing. Distance metrics are nice to have and seem to have the potential for creating a table-wise metric [220], so presenting them is important. Second, for privacy evaluation, we believe that Differential Privacy in itself is not a guarantee of protection for real patients. More research and depth should be employed when presenting results for such generators; record-linkage and attribute disclosure can provide extra guarantees. Thirdly, a clinical evaluation should be done as well to understand if the synthetic patients are a reality in the clinical setting. Since the correlations could be correct but clinically (or biologically) they might not make sense. Finally, in the scope of this paper, only GANs were assessed, but there are more mechanisms for generating data, and could be interesting to assess how all of them perform on the same datasets. There are other methods for handling the mixed data types that regularly appear in clinical settings, like Variational Autoencoders (VAEs) Gaussian Mixtures, BN, and imputation mechanisms, making them excellent candidates for this assessment.

### 3.1.6 Conclusion

In this paper, we had the opportunity to survey the current framework for generating tabular data using GANs and which ones were already tested in the healthcare setting. We summarised the utility and privacy metrics employed, and the datasets used to measure them. We analyzed the code availability and made suggestions for further work on cataloging, comparing, and assessing synthetic health data generators. A survey with a global benchmark of methodologies, despite being arduous, could yield great results for the community and take the aim of this paper further.

## 3.2 Pulling the current metrics of assessing datasets

This section is based on the paper entitled "Dataset Comparison Tool: Utility and Privacy". This work followed the work on section 3.1, where we compiled ways of assessing the utility of syn-

thetic data. We understood that the mechanisms were far from consensual and a tool could be of use to merge all of this into a single file and report about data. Our purpose was to facilitate health data owners and legal responsible to understand how similar and protective a dataset was regarding the original one.

### 3.2.1 Introduction

Synthetic data can be defined as data that has no connection with a real-world phenomenon or event. It did not originate from a process in the real world, but rather a synthetic one. The idea is that synthetic data can have similar properties with real data, without needing to have an independent process for its generation. Synthetic data has been used over the years for several usages, but in healthcare is still not very used. However, this scenario seems to be changing. It can be used for several use cases namely [94]; i) Software testing, ii) educational purposes, iii) ML, iv) regulatory, v) retention, vi) secondary and vii) enhanced privacy.

Software testing relates to using synthetic data to create use cases for software testing. This can be used for the development or pre-production stages for example. Often the data needed is not available on-demand and a synthetic generator of reliable data could be useful. Educational purposes relate to, at least, two different scenarios. One is for onboarding of employees [94], the other is related to healthcare students for using health information systems and creating mechanisms for providing reliable data on-demand. ML is one of the areas where synthetic data has more widespread usage, where data augmentation through data synthesis is already common. It can be used for class imbalance, sample-size boosting, or machine-learning algorithm testing. Regulatory purposes could be important as well, with the rise of AI as medical device systems and synthetic data could be used to properly evaluate these systems under controlled environments. Retention can be an important case for synthetic data as well, since personal data must not be kept more than it would be necessary. Synthetic data generators can be of use, where the original data can be deleted and a generator kept for further usage, given that the privacy mechanisms are properly employed. Secondary uses relate to using synthetic data to share data with academia or industry. Simulacrum [85] is a nice example of how the NHS uses these mechanisms to help scientists get a better grasp of data before having to fill in documentation to query the real data. The same occurs for Integraal Kankercentrum Nederland (IKNL), which has a synthetic version of the cancer registry for scientific purposes [93] and the Healthcare Products Regulatory Agency (MHRA) that uses synthetic data as well for its CPRD real-world evidence [58].

Finally, an aspect that is underlying all these applications is the promise that synthetic data can be used to improve privacy. Even though specially tweaked data generators can be used to create more privacy-aware datasets, it will be inherently always at the cost of some utility [178]. So, even though synthetic data is not the silver bullet as primarily thought, synthetic data generation can be undeniably used to help create more private data for all the use cases seen above, at the cost of its utility. As for proper methods of evaluating security and utility, there are, for now, open research questions. At the present time, it is still complicated to properly assess the utility of the generated data. We have qualitative and quantitative methods. Qualitative methods are related to plots, and

quantitative are related to some value that defines an evaluation metric. These quantitative metrics can be applied to equal columns from each data set, pair of columns from each dataset or applied over the whole dataset. As for privacy metrics, the metrics rely on duplicates. Full duplicates or membership inference related.

So in this paper, we developed a data pipeline for data analysis in order to create a report for providing several metrics for data utility and privacy.

### 3.2.2 Methods

The pipeline relies on Python and latex for creating the document. It relies also on several packages that implemented methods for evaluating data, namely *scipy* [199], *sdmetrics* [143] and *scikit-learn* [146] and *mlxtend* [159]. Its basis is related to uploading 2 datasets, and a report in pdf is produced. Being that is based on programmatic code, it can be easily converted into Application Programming Interface (API). The report has a section for dataset description, columns removed due to high-null, and a brief variable overview. Then a null comparison is done by column and dataset. Following this is the utility subsection. Firstly by visual methodologies: heat maps for the correlation, bar plots for categorical, density plots for continuous, and a pair plot for an overview. As for the quantitative utility evaluation, we divided it column-wise, pair-wise, and table-wise. The first comprehends the *Kolmogorov–Smirnov* test for continuous and  $\chi^2$  test for categorical variables. Distance metrics were also applied to categorical columns. First, they are transformed into distributions and then distance metrics are applied. The results is a descriptive overview of the distance metrics, having minimum value, average, max value, and standard deviation. The distance metrics selected were *Jensen-Shannon Divergence*, *Wasserstein distance*, *Kullback–Leibler divergence*, and entropy. As for pair-wise metrics, we used a discrete and continuous *Kullback–Leibler divergence*. In this, two pairs of continuous columns are compared using *Kullback–Leibler divergence*. For this, they are put into bins for further application. The same is applied to categorical columns without binning. As for table-wise metrics, first, we used likelihood metrics. We fitted several Gaussian Mixture models or BN models to the real data and then calculated the likelihood of the synthetic data belonging to the same distribution. The metrics are likelihood for the Gaussian mixture and Bayesian models and log-likelihood for the Bayesian model as well.

Then we used machine-learning models (linear regression and decision trees) to assess how similar models behave on both datasets. First, we tested on the same dataset in order to compare evaluation metrics. Then we cross-tested, meaning that the training set was drawn from one dataset and the test set was drawn from the second dataset. Finally, data privacy constraints duplicate evaluation, duplicate existence by removal of a single column and a record linkage approach. With the record linkage, we define a record linkage blocking ("age" in the example) and then try to match rows from the synthetic dataset to the real, with varying known attributes. Then matrix, euclidean and cosine distance was also calculated. Even though it is used for privacy evaluation, by definition, we could also use it for utility assessment. For proper assessment, the continuous and categorical variables should be indicated at the start of the code. The metrics are listed in the table 3.3.



Table 3.3: Metrics Assessed

Metric	Method	Context
Bar Plot	visual	utility
KDE Plot	visual	utility
Heat-map	visual	utility
Pair-plot	visual	utility
KS test	column-quantitative	utility
ChiSquared Test	column-quantitative	utility
Kullback–Leibler divergence	column-quantitative	utility
Jensen-Shannon Divergence	column-quantitative	utility
Wasserstein distance	column-quantitative	utility
Entropy	column-quantitative	utility
DiscKLD	table-quantitative	utility
ContinuousKLD	table-quantitative	utility
BNLikelihood	table-quantitative	utility
BNLogLikelihood	table-quantitative	utility
GMLogLikelihood	table-quantitative	utility
Same dataset ratio	table-quantitative	utility
Support rules	table-quantitative	utility
Different dataset validation	table-quantitative	utility
Duplicates	quantitative	privacy
Duplicate minus 1	quantitative	privacy
Record Linkage	quantitative	privacy
Matrix distance	quantitative	privacy/utility
Cosine distance	quantitative	privacy/utility
Euclidean distance	quantitative	privacy/utility

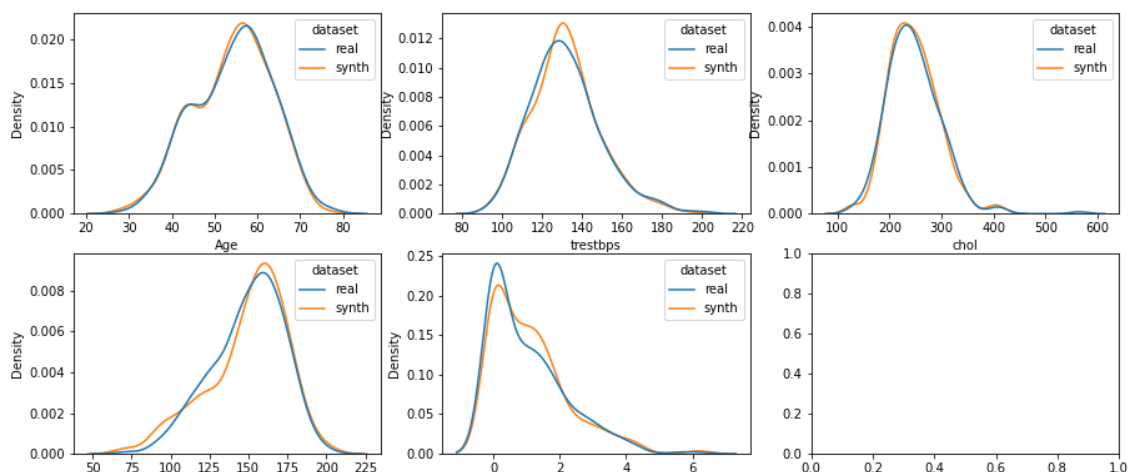


Figure 3.3: Continuous Variables plotted

### 3.2.3 Results

A trial example of comparing data is available for data in the UC Irvine Machine Learning Repository (UCI) repository, namely the heart disease dataset [95]. The synthetic data was created by using the synthpop package [135]. The variables evaluated are listed in table below. The code can be seen in <https://github.com/joofio/dataset-comparasion-report>. As an example. The image for visual analysis for categorical (figure 3.2) and continuous variables (figure 3.3).

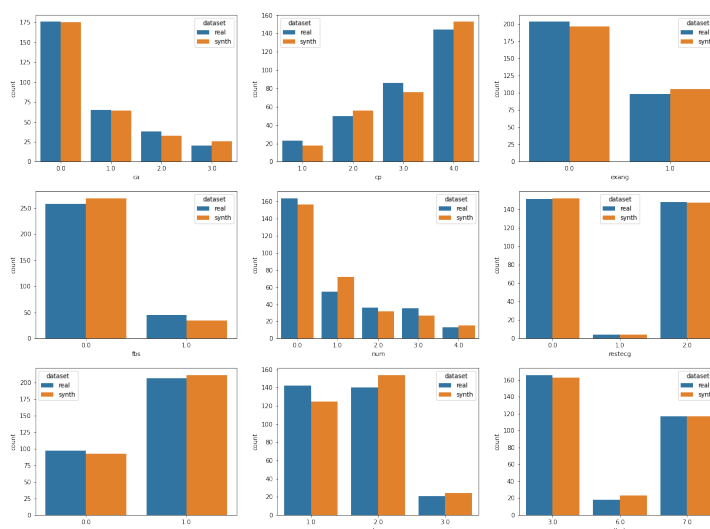


Figure 3.2: Categorical Variables plotted

### 3.2.4 Discussion & Conclusion

The data possible to create to evaluate similarities between two datasets is important not only for synthetic vs real datasets. For example in distributed learning, where different silos exist, with similar or even equal features, a method for evaluating the similarities can be useful for understanding how the populations are similar between them, trying to shed light on the most similarities among them, or different in order to understand the differences in the silos or data acquisition inside them. Furthermore, the differences can be assessed on a more granular level. The column-wise similarities can be different from the inter-column similarities and this in itself, can be a metric of interest regarding the quality of the synthetic data and its generator.

With this work, we hope to help institutions and academics get access to a benchmark of the datasets provided in order to leverage synthetic data in the healthcare space. Finally, we hope this work helps other researchers reach an evaluation metric that could be a unique and clear response to the question of how similar two datasets are.

## 3.3 Can we use machine learning feature to compare datasets?

This section is based on the paper entitled "Using Machine Learning Models' feature importance to assess dataset similarity". The reasoning behind this paper was the results of 3.1, where we felt that evaluation metrics for synthetic data could be improved. Better yet, we felt that the comparison of two datasets (that shared the same columns) could be done in a more robust way. Being that the current gold-standard was cross-validation which was not bound to any number range and the significance of the result could not be easily interpretable. We used the feature importance of several ML models to compare datasets and concluded that it was a valid alternative to the traditional metrics.

### 3.3.1 Introduction

In recent years, the use of AI and ML algorithms has gained increasing prominence in healthcare research and practice. One of the key requirements for the successful application of these methods is access to large, high-quality datasets. However, in many cases, the availability of such datasets can be limited due to issues around data privacy, security, and ethical concerns [43]. To address this challenge, synthetic data has emerged as a promising solution. Synthetic data refers to artificially generated data that closely mimic the statistical properties and patterns of real-world data [130].

Synthetic data has the potential to overcome many of the limitations associated with real-world data, such as the lack of sufficient data volume, noise, and privacy concerns. Even though there are still doubts if the privacy part is the silver bullet sometimes referred to [179], the upsampling part is a standard use for years now. However, the quality of synthetic data generated by various techniques can vary significantly, and it is essential to assess the quality of synthetic data before its usage. In healthcare, the assessment of synthetic data is crucial to ensure that it can provide valid insights and inform decision-making processes.

The assessment of synthetic data in healthcare is essential for its successful use in various applications, such as developing predictive models, testing algorithms, and conducting clinical trials. The use of synthetic data can significantly enhance the efficiency and effectiveness of healthcare research and practice. However, it is crucial to ensure that the synthetic data used in these applications are of high quality and validated to provide reliable and valid insights. The evaluation of synthetic data quality involves comparing its statistical properties and patterns with those of the original data. We can assess how similar columns are to each other through several statistical tests, and then we can infer some inter-column properties with methods like cross-validation, where two datasets are split into train tests and cross-tested and then the ratio between the evaluation result of both datasets is used as a metric [130, 77]. However, this methodology is a big proxy for such an inter-column relationship. Can we try to provide a better metric than this one to evaluate how similar are the inter-column relationship of two distinct datasets? In this paper, we suggest using feature importance values to create a more explainable and reasonable metric for inter-column relationships.

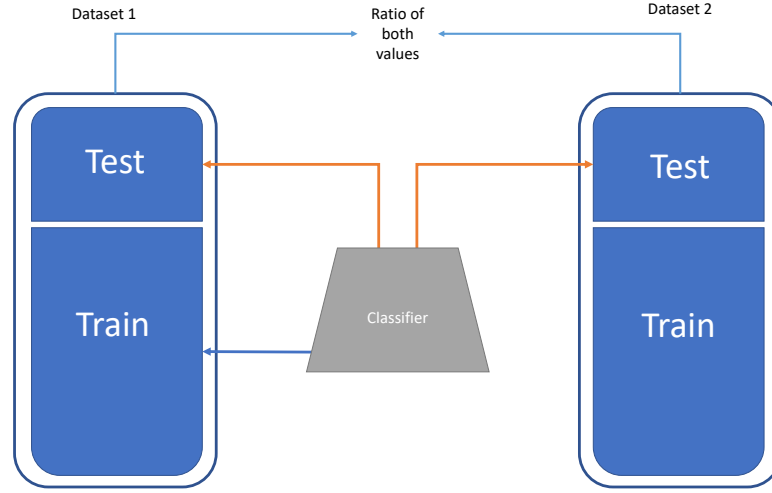
### 3.3.2 Rationale and Related Work

Recently there has been a series of works related to assessing how synthetic data generators behave with data like the work of Emam et. al [67] that especially focused on utility metrics for synthetic data generators. At the moment, comparing data is based on intra-columns and inter-columns relationship. The intracolumn relationship is assumed as something that compares equal columns between datasets, with highly known statistical methods like chi-squared or Kolmogorov Smirnov like done in the works of [49] among many others, acting more like sanity checks than anything else. Other known metrics are distance-based metrics like Jensen Shannon divergence, Wasserstein Distance, Bhattacharyya Distance or Hellinger distance, which are based on the calculation of the distance between distributions like seen in the works of several teams [222, 46, 20].

However, regarding inter-column relationships, the metrics applied are often very different across papers. One example of trying to capture inter-column relationship is about the use of propensity score [167, 130] where a classifier is trained to the merged datasets, with the added variable of the original dataset (i.e., 1 for real and 0 for synthetic). The model is trained and the propensity Mean square error is the mean squared difference of the estimated probability from the average prediction. Most recently, a unified metric appeared as the sum of other metrics known as described in the work of Chundawat et. al, [47], known as TabSynDex. Other examples are likelihood of fitness like in the works of [218], coverage support [77] or very specific metrics implemented for evaluating specific data generators. However, the most used metric is cross-validation, which takes two datasets, one that is real and a second which is synthetic and we split both into train and test and train a machine learning model on the real data training set, then we test the model on both test sets. Then a ratio is created, rendering the actual value. This methodology, even if gold-standard at the moment for this type of study, has some liabilities since this value can be a bit erratic, and even above one since the evaluation metric could be better on the second dataset and we don't have a clear grasp of what that can represent in terms of dataset similarity.

The image 3.4 represents this in detail. Several works used this metric as the comparing metric [130].

Figure 3.4: Cross-Validation of datasets



### 3.3.3 Materials & Methods

#### 3.3.3.1 Method Overview

For this work, our goal is to test several metrics based on the ranking of feature importance of a trained model. Normalized Discounted Cumulative Gain (NDCG) [206] which is the sum of the true scores ranked in the order induced by the predicted scores, after applying a logarithmic discount. Then divide by the best possible score to obtain a score between 0 and 1. It is calculated by

$$\text{NDCG} = \frac{\text{DCG}(P)}{\text{IDCG}(P)} \quad (3.3)$$

where  $\text{DCG}(P)$  is the Discounted Cumulative Gain and  $\text{IDCG}(P)$  is the Ideal Discounted Cumulative Gain.

Cohen's kappa coefficient [48] is a statistic that is commonly used to assess the level of agreement between two or more raters or evaluators who are providing categorical ratings or rankings of a set of items. So, we want to use to assess if it could be of use to check how similar the ranking of the features is, using the numbers as categorical.

$$\kappa = \frac{P_o - P_e}{1 - P_e} \quad (3.4)$$

where  $P_o$  is the observed agreement between the two raters and  $P_e$  is the expected agreement between the two raters by chance.

We also intend to use the  $R^2$  to check if the explainability changes across datasets.

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (3.5)$$

where  $y_i$  are the observed values of the dependent variable,  $\hat{y}_i$  are the predicted values of the dependent variable,  $\bar{y}$  is the mean of the observed values of the dependent variable and  $n$  is the number of data points.

Then we intend to use ranking metrics, namely Kendall tau, weighted Kendall tau and RBO. Kendall tau is a measure of correlation that measures the similarity between two rankings. It is commonly used in statistics and data analysis to evaluate the agreement or disagreement between two sets of rankings.

The Kendall tau coefficient [105] is defined as the difference between the number of concordant and discordant pairs of observations, divided by the total number of pairs. A concordant pair is a pair of observations that have the same ranking order in both sets, while a discordant pair is a pair of observations that have opposite ranking orders. The Kendall tau coefficient ranges from -1 to 1, where -1 represents perfect negative correlation, 0 represents no correlation, and 1 represents perfect positive correlation.

$$\tau = \frac{\text{number of concordant pairs} - \text{number of discordant pairs}}{\text{total number of pairs}} \quad (3.6)$$

Weighted Kendall tau [197] is an extension of Kendall tau that takes into account the importance or weight of each observation in the rankings. In some cases, some observations may be more important than others, and their positions in the ranking may have a greater impact on the overall correlation. Weighted Kendall tau assigns a weight to each observation, and the correlation is calculated based on the weighted concordant and discordant pairs.

$$\tau_w = \frac{\sum_{i < j} w_{ij} \cdot \text{sgn}(x_i - x_j)}{\sum_{i < j} w_{ij}} \quad (3.7)$$

where  $w_{ij}$  is the weight associated with the pair  $(x_i, x_j)$  and  $\text{sgn}(\cdot)$  is the sign function. Rank-biased overlap (RBO) [208] is a measure of similarity between two ranked lists or rankings. It takes into account the order of items in the two lists, and it can be used to evaluate the quality of search results, recommendations, or any other kind of ranked list it has been shown to be more robust and accurate than other similarities measures such as Kendall tau or Spearman's rank correlation coefficient.

$$\text{RBO} = (1 - \rho) \cdot \sum_{d=1}^{\infty} \left( \frac{g_d}{d} \right) \cdot \rho^d \quad (3.8)$$

where  $\rho$  is the weight,  $g_d$  is the gain at depth  $d$ , and  $\sum_{d=1}^{\infty}$  indicates the summation over all depths.

Finally, we intend to use text-distance metrics. The theory behind this experiment is to treat the ordered columns in a ranking manner and apply text-distance metrics to check the distance be-

tween the two. Levenshtein distance [131] is the minimum number of single-character insertions, deletions, or substitutions required to transform one string into another.

Damerau-Levenshtein distance [131] is similar to Levenshtein distance but also includes the transposition of two adjacent characters as an allowable operation.

The hamming distance [83] is a measure of the difference between two strings of equal length, defined as the number of positions at which the corresponding symbols are different.

Jaro-Winkler distance [131] is a string similarity measure that takes into account the number of matching characters, the number of transpositions, and the length of common prefixes, with a higher weight given to the common prefix.

```

for i in number of columns to test do
  for rep in 10 repetitions do
    permutate values in i columns
    for dataset in dataset pair do
      for target in dataset columns do
        Train-Test Split (95:5) model fit to train get feature importance per column
        Create an ordered rank of features
      end
    end
  end
end

```

**Algorithm 1:** Testing similarity scores in tabular datasets

The algorithms chosen were decision trees, random forests, SVM, KNN, and linear regression/logistic regression as implemented in the scikit-learn package [146]. The text distance metrics were implemented by the text-distance package [139]. Kendall tau, weighted Kendall tau were used as implemented by scipy [200] and RBO, as implemented in [41].

### 3.3.3.2 Data used

We used 5 datasets from the UCI repository. The ones chosen were related to healthcare and were heart disease [95], thyroid disease [155], liver disorders [72], breast cancer [212] and the primary tumour dataset [224]. We made minimal preprocessing on the datasets, namely removing the missing variables by imputing the mean on continuous variables and mode on categorical. We also created a synthetic dataset by applying the synthpop package to this data [135]. With this package, all variables were synthesised using the "cart" method, which is rpart implementation of a CART model.

### 3.3.4 Results

With the method described in the algorithm 1, we created a figure where the metrics are presented with increasingly different datasets: figure 3.5.

The number of repetitions and how that impacts the variance of the scores is shown in figure 3.7.

Figure 3.5: Plot showing the decrease of the metric over increasingly changed datasets. The X axis represents the number of columns mutated. The Y axis represents the value of the metric and the hue represents the algorithm used to calculate the metric.

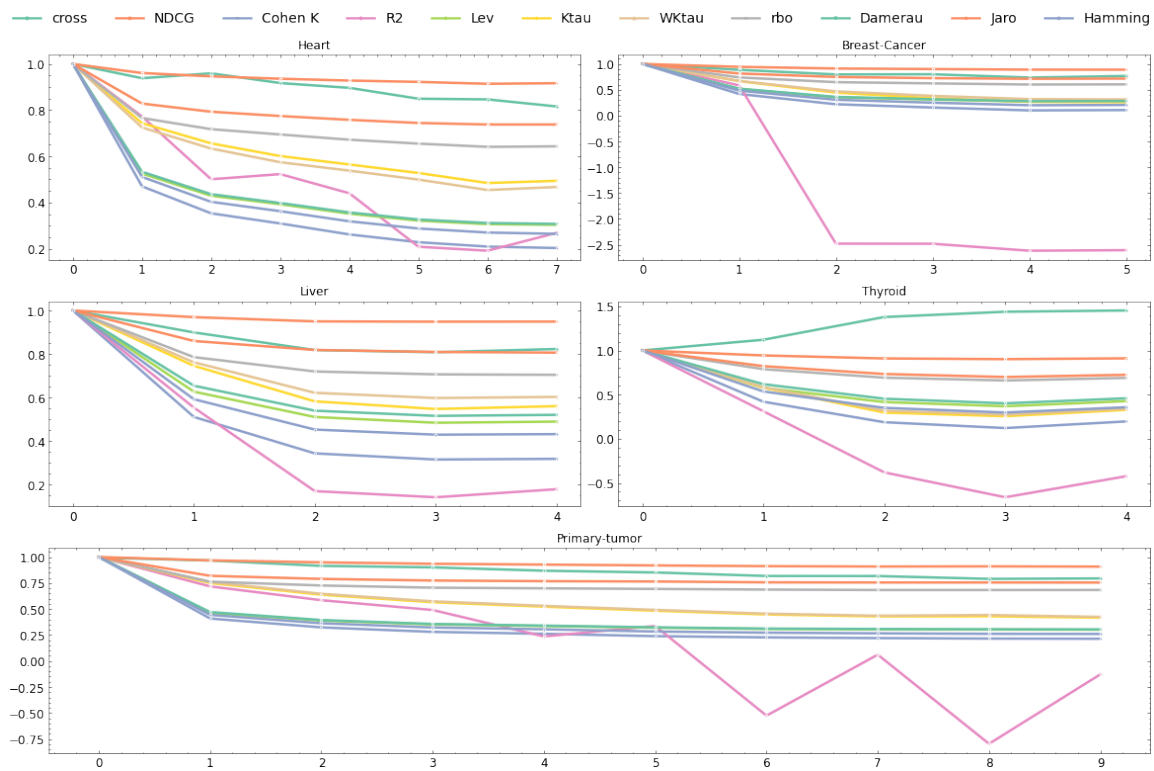




Figure 3.6: Plot showing the values at 50% columns mutated across all datasets and algorithms per metric type

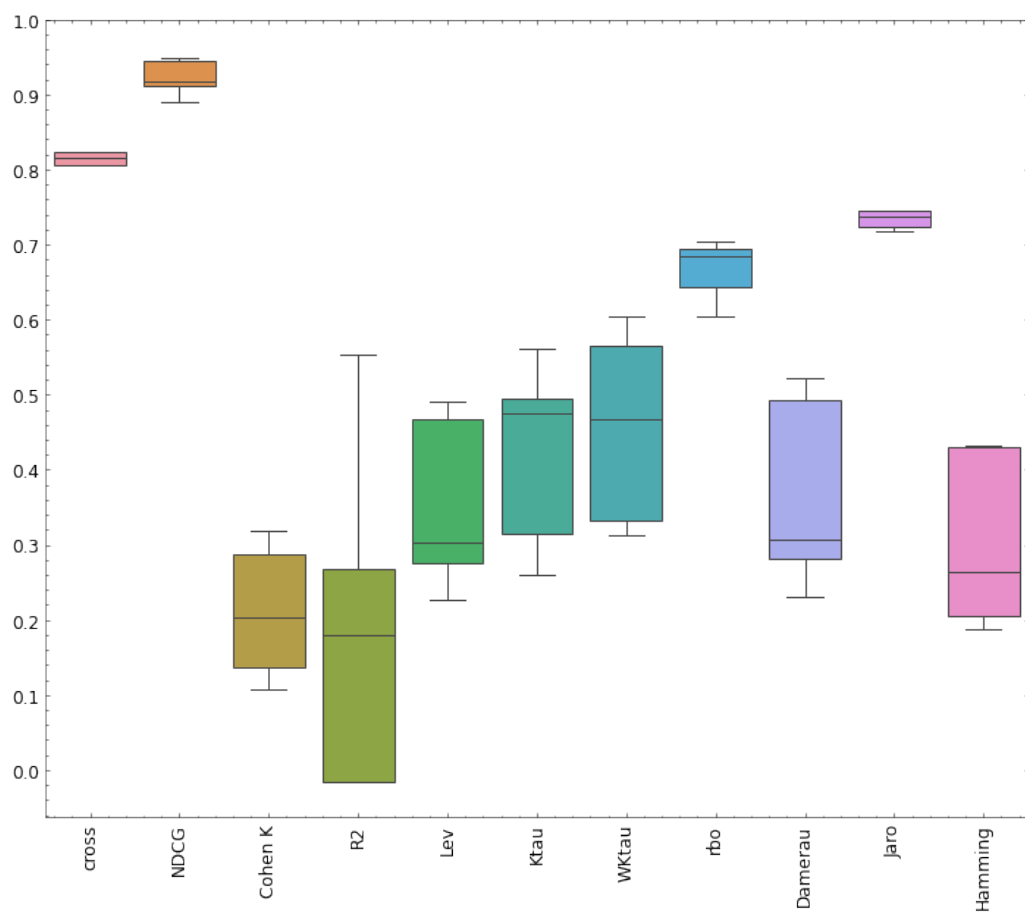
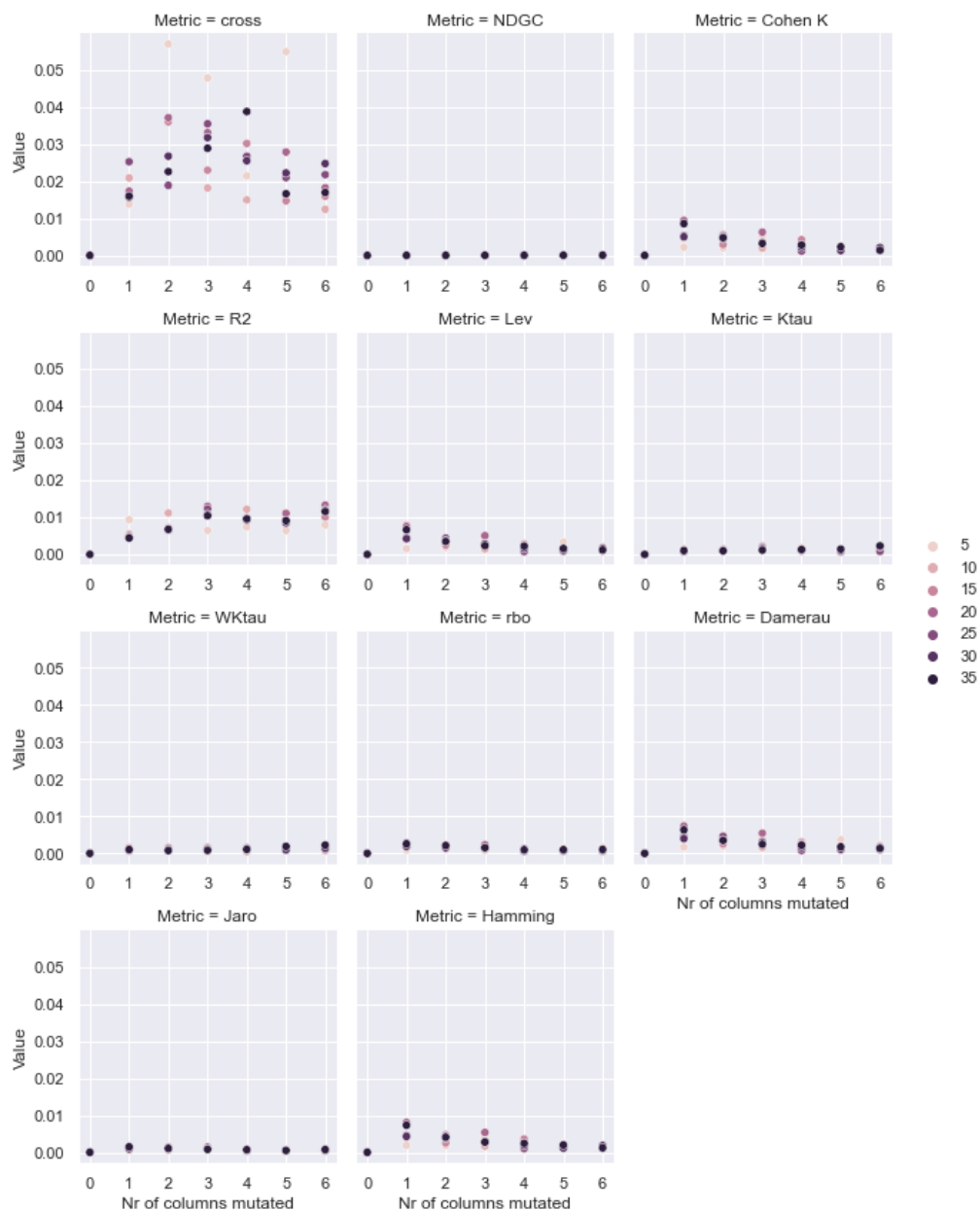
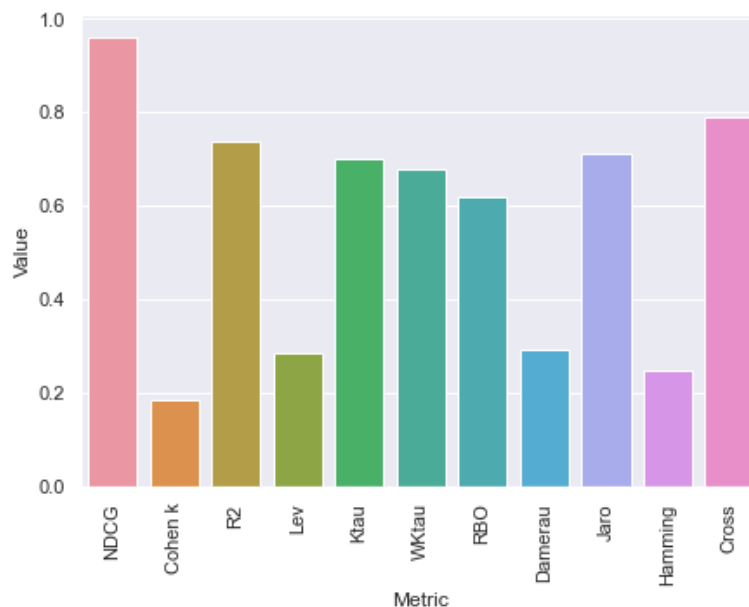


Figure 3.7: Plot the variance of different repetitions for every metric and the number of different columns changed. X is the number of columns mutated. Colour is the number of repetitions for each mutation and Y is the variance of the data.



As for the test for the synthetic and real datasets, the results are displayed in figure 3.8.

Figure 3.8: Result on a synthetic and real data



### 3.3.5 Discussion

With the results found, we feel that are better alternatives to cross-validation. At least Kendall tau, Weighted Kendall tau, and RBO seem like alternatives to cross-validation. Firstly, they seem to be directly connected to a difference in the dataset, secondly, they are a 0-1 metric and thirdly, variance across different iterations is also lower. From these, the metrics based on ranking metrics seem to work best, where Kendall tau, weighted Kendall tau and RBO have better performance than the rest. As for the variance with the number of repetitions, we also see that the ranking-based metrics have good stability, while the cross-validation and text-based metric have higher variability with a low number of repetitions (even if low) - figure 3.7. Text metrics also have a suitable performance, even though they have a drastic drop with only one column mutated (figure 3.5).

### 3.3.6 Conclusion

Comparing two tabular datasets has been growing in demand in the past year mainly because of the increase in popularity of tabular data synthesis methods which have exhibited the potential in generating valuable synthetic data. However, due to the absence of a uniform metric, evaluating different methods has been inconsistent. This research proposes some alternatives for assessing synthetic tabular data's utility. Ranking metrics (Kendall Tau, Weighted Kendall Tau, and RBO) and  $R^2$  have shown the potential to capture inter-column relationships in a more consistent way than cross-validation. They could become a useful tool for comparing statistical methods of generating synthetic tabular data. Furthermore, this metric can aid in evaluating these generators'

training, providing insights into improving synthetic data quality. The proposed metrics open up possibilities for future research to enhance tabular data synthesis methods and compare two datasets overall.

### 3.4 Data quality Metrics

This section is based on the paper entitled "so and so"

#### 3.4.1 Introduction

With the wide spreading of healthcare information systems across all contexts of healthcare practice, the production of health-related data has followed this incremental behaviour. The potential for using this data to create new clinical knowledge and push medicine further is tempting [123]. However, to correctly use the data stored in EHRs, the quality of the data must be robust enough to sustain the clinical decisions made based on this data. The issue is that data quality is not a straight line and is very context aware. The threshold and dimensions required to classify the quality of the data depend on the purpose that we intend to use that very same data [202]. These uses can be very distinct and have different impacts as well. For one, we can use data to support day-to-day decisions regarding individual patients' care [195]. These decisions can include ones based on recorded information to understand a patient's history, clinical decision support systems based on this data, or even using the data to help support a more macro, public health-oriented decision. Another area is using information for management purposes. The data can be used by management bodies and regulatory authorities to extract metrics regarding the quality of care or reimbursement purposes. Thirdly, data can be used for research purposes, namely observational studies and, more recently, to support clinical trials through real-world evidence analysis [54, 195, 211]. So, all the EHR data-based decisions can only be as good as the data supporting them. Several studies have already warned about the lack of data quality in EHRs and how this can be a significant hurdle to an accurate representation of the population and potentially lead to erroneous healthcare decisions [163, 99, 92, 223, 107, 73].

There are several steps in the data lifecycle that can be prone to error, from data generation, where the data is registered by healthcare professionals, passing by data processing, whether inside healthcare institutions or by software engineers aiming to reuse data, to data interpretation and reuse, where investigators try to interpret the meaning of registered data [211]. So, with all of the data's possible uses added to the several steps that can introduce errors throughout the data lifecycle, data quality frameworks and sequential implementations can have very distinct approaches and methodologies to assess data quality. Data quality tools for checking data being registered live to support day-to-day decisions will be significantly different from one whose only purpose is to provide quality checks for research purposes. So, methodologies to tackle these issues are necessary for guaranteeing the quality of healthcare practice and the knowledge derived from EHR data. Consequently, in this paper, we propose:

- Create a tool for identifying data quality issues in obstetrics EHRs;

- Enlighten on the issues that can appear with a full deployment of such a tool
- Suggestion of a creation of a single score for data quality for comparison of high-quality and low-quality records in a database.
- Assess how such a tool can work in early-stage real-world scenarios and how to work with obstetricians to improve data quality.
- Identify data quality issues on obstetrics data

### 3.4.2 Background and Related Work

There is already a significant number of papers trying to define data quality assessment frameworks for EHR data, all of them plausible and recommendable, already described in other papers [28]. The literature has over 20 different methods, descriptions, and summaries of different frameworks over the years. Some may be highlighted from the review from Weiskopf et. al, [210], where five data quality concepts were identified over 230 papers: Completeness, Correctness, Concordance, Plausibility and Currency. Then Khan et al. tried to harmonize data quality assessment frameworks, which simplified all previous concepts into three main categories: Conformance, Completeness and Plausibility and two assessment contexts: Verification and Validation [100]. Then a review of Bian et al. [28] expanded on the previous ones, categorizing data quality into 14 dimensions and mapping them to the previous most known definitions. These were: currency, correctness, plausibility, completeness, concordance, comparability, conformance, flexibility, relevance, usability, security, information loss, consistency and interpretability. Despite all of these comprehensive works, there is still no consensus regarding which one is best or which has taken the lead in usage. Moreover, looking at all of the descriptions related in the literature, a significant portion of concepts are overlapping and sometimes hard to conceptualize such dimensions in practice.

As for implementations, there are already some available, such as the work from [148] where a tool created by primary care in the Flanders was built to assess completeness and percentage of values within the normal range. The work from Liaw et al. [113] already reviewed some data quality assessment tools, like tools from OHDSI [89] or TAQIH [11]. Additionally, we found some others with similar purposes and characteristics like the work presented dataquaieR [173], an R language-based package that can assess several data quality dimensions in observational health research data. Also, the work from Razzaghi et al. developed a methodology for assessing data quality in clinical data [162], taking into account the semantics of data and their meanings within their context. Furthermore, the work from Rajan et al. [156] presented a tool that can assess data quality and characterize health data repositories. Parallel to this, Kaspner et al. created a tool called DQASStats that enables the profiling and quality assessment of the MIRACUM database, being possible to integrate into other databases as well [102]. However, these tools are not meant to be used at the production level, assessing data as it is being registered or outputs reports for human consumption and not a quantitative metric for metric comparison. Furthermore, none of the non-agnostic tools were designed for obstetric EHR data. Finally, we have not seen, until the

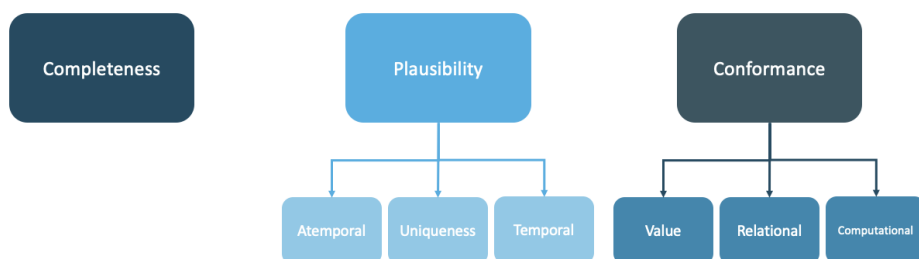
moment of this paper, any implementation that used machine learning to evaluate the correctness of the value.

### 3.4.3 Materials

The data was gathered from 9 different Portuguese hospitals regarding obstetric information: data from the mother, several data points about the fetus and delivery mode. The data is from 2019 to 2020. The software for collecting data was the same in every institution, and the columns were the same, even though the version of each software differed across hospitals. Across the different hospitals, data rows ranged from 2364 to 18177. The sum of all rows rendered 73351 rows. The data dictionary is in appendix A.1.

For this purpose, we took the Khan harmonized framework since we understood it as simpler to communicate we feel that the three main categories are indeed non-reducible, which makes sense from an organizational standpoint. Furthermore, the work done by Khan et al. with mapping to already existing frameworks could help compare this work with others that felt the need to use other frameworks. With this in mind, we will use three main categories, Completeness, Conformance and Plausibility. Completeness relates to missing data. Conformance relates to the compliance of the data representation, like formatting, computational conformance and other data standards implemented. Plausibility relates to how believable the values are.

Figure 3.9: Dimensions of data quality



### 3.4.4 Methods

We wrote all the code in Python 3.10.6 with the usage of the scikit-learn library for preprocessing, and evaluation [146]. For plausibility, a Bayesian network was used. We used this model due to the possibility of using a single model for classifying the plausibility of all columns and due to its interpretable nature. The networks were created with the pgmpy package [13]. We also added the outlier-tree method [?] which tries to integrate a decision tree that "predicts" the values of each column based on the values of each other column. In the process, every time separation is evaluated, it takes observations from each branch as a homogeneous cluster to search for outliers in the predicted 1-d distribution of the column. Outliers are determined according to confidence intervals in this 1-d distribution and need to have large gaps in order to be marked as outliers in the next observation. Because it looks for outliers in the branch of the decision tree, it knows the

conditions that make it a rare observation relative to other observation types corresponding to the same conditions, and these conditions are always related to target variables (as predicted by them). As such, it can only detect outliers described by decision tree logic, and unlike other methods such as isolation forests, it can not assign outlier points to each observation, or detect outliers that are generally rare, but will always provide human-readable justification when it recognizes outliers.

As for preprocessing, all null representations were standardized, we also removed features with high missing rates ( $> 80\%$ ). The imputation process was performed with the median for continuous and a new category (NULLIMP) for categorical variables.

For the usage of the Bayesian network in particular, the continuous variables were discretized into three bins defined by quantile. We defined three as the number of bins in order to reduce the number of states in each node of the network. The evaluation was done with cross-validation with 10 splits and two repetitions for each column as the target.

As for Z-Scores, they were defined for all continuous variables based on the interquartile range. Then, rows were also assessed with distance analysis, with Local Outlier Factor and Elliptic Envelope from scikit-learn and the outlier-tree algorithm. We also added a rule engine, using the *great\_expectations* package. Rules were defined by the team, focusing on impossible numbers present in age, weight, or relationship between variables. As for missing information was created with all the data, creating the scoring based on the inverse of the missing percentage. Missing detection was based on primary key variables. For completeness, we used the inverse of the percentage of nulls in the training set. The API for serving the prediction models was developed with FastAPI. So, the methods applied in terms of the DQA framework shown in figure 3.9 are described in the table 3.4.

Table 3.4: Implemented Methods

Category	Subcategory	Method
Completeness	N/A	Score by the inverse percentage of missing in the train data
Plausibility	Atemporal Plausibility	Bayesian model prediction based on the other values of row
Plausibility	Atemporal Plausibility	Z-score for column value based on IQR train data
Plausibility	Atemporal Plausibility	Elliptic Envelope
Plausibility	Atemporal Plausibility	Local Outlier Factor
Conformance	Value Conformance	Manual Rule engine
Plausibility	Atemporal Plausibility	Manual Rule engine
Plausibility	Atemporal Plausibility	outlier-tree

The method of scoring was to obtain a single value that could grasp the quality of the row or patient. To assess the tool's usefulness, we will implement it in a production environment and

collect metrics regarding the data being produced. Then we intended to present some results to selected obstetrics clinicians for them to assess how likely the information is to be suitable for usage. We will also compare the results with the ones from the model to make sanity checks regarding the model's performance and adequacy. We aim to use Kendal Tau and Average Spearman's Rank Correlation Coefficient. Kendall Tau is a non-parametric statistic used to measure the strength and direction of the association between two ordinal variables. It calculates the difference between the number of concordant and discordant pairs of observations, normalized to ensure a value between -1 (perfect disagreement) and 1 (perfect agreement). Spearman's rank correlation coefficient is a non-parametric measure that assesses the strength and direction of a monotonic relationship between two ranked variables. It is based on the ranked values of the variables rather than their raw data, producing a value between -1 (perfect inverse relationship) and 1 (perfect direct relationship).

### 3.4.5 Results

A Bayesian network with structure and parameters learned from the training dataset reached an average Area Under the Receiver Operating Characteristic Curve of 0.857. The results are in the table 3.5.

Table 3.5: Validation Results: Column acronym with AUROC along with 95% CI

AP	0.944	[0.943, 0.945]	VNH	0.894	[0.893, 0.895]
AG	0.797	[0.778, 0.816]	TPEE	0.816	[0.815, 0.816]
EA	0.969	[0.968, 0.969]	AA	0.751	[0.743, 0.758]
CA	0.958	[0.958, 0.958]	GR	0.931	[0.93, 0.932]
IA	0.638	[0.637, 0.638]	V	0.983	[0.982, 0.983]
PI	0.881	[0.88, 0.881]	TP	0.866	[0.865, 0.868]
IMC	0.881	[0.881, 0.882]	VCS	0.79	[0.789, 0.791]
NRC	0.75	[0.75, 0.75]	ANP	0.942	[0.938, 0.946]
IGA	0.968	[0.968, 0.969]	GS	0.514	[0.507, 0.52]
SGP	0.974	[0.974, 0.974]	S	0.896	[0.896, 0.897]
VA	0.974	[0.974, 0.974]	VP	0.771	[0.77, 0.772]
TG	0.728	[0.726, 0.73]	TPNP	0.952	[0.951, 0.952]
<b>Average 0.857 [0.846, 0.868]</b>					

The network is as represented in figure 3.10.

As for the rules created, they were conformance-based, like the format of dates, and conformance to the value set (i.e. Robson group, bishop scores, or delivery types). We also added plausibility rules, like expected values for BMI, weight, and gestational age. We also added plausibility for the relationship between columns, namely weight across different weeks of gestation. We have also added a relationship of greatness between ultrasound weights more than 5 weeks apart. The method of calculating the final score is stated in figure 3.11.



Figure 3.10: Network learned

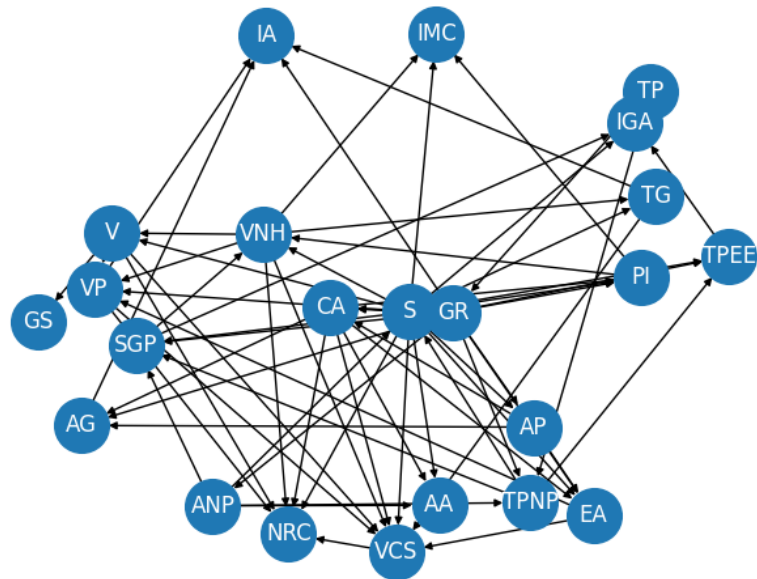
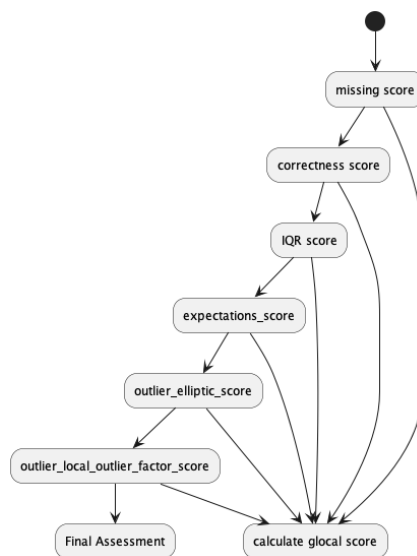


Figure 3.11: Workflow for creating the final score and which elements are used to do so.



### 3.4.6 Deployment & Validation

The purpose of this model is to serve as an API for usage within a healthcare institution and act as a supplementary decision support tool for obstetrics teams. Although a concrete, vendor-specific information model and health information system were initially used, our goal is to develop a more universal clinical decision support system. This system should be usable across all systems involved in birth and obstetrics departments. Therefore, we constructed it using the Health Level 7 (HL7) Fast Healthcare Interoperable Resources (FHIR) R5 version standard. This approach simplifies the process of API interaction. Rather than utilizing a proprietary model for the data, we based our decision on the use of FHIR resources: Bundle and Observation. These resources handle the request and response through a customized operation named "\$quality\_check". We intend to publish the profiles of these objects to streamline API access via standardized mechanisms and data models. The model then makes use of the customized operation and of several base resources to construct a FHIR message, which are: Bundle, MessageHeader, Observation, Device. Observation is where the information about the record is contained, Device contains information about the model, and MessageHeader is used to add information about the request. Finally, the Bundle is used to group all of these resources together. The current version of the profiles can be accessed at this URL: <https://joofio.github.io/obs-cdss-fhir/>.

For validation, we deployed the tool in docker format in a hospital to gather new data. We gathered 3231 new cases and returned a score for quality as exemplified in figure 3.12. Being that the score is from 0 to 1, the average score was 0.23 and IQR was 0.03. We also used the clinician from one of the hospitals that we get data from and asked this clinician to assess 10 records in terms of quality. We gathered the 10 records at random and asked the clinician to assess them in terms of quality. Our purpose was then to compare the rankings of each evaluator; the model and the clinician, to assess how similar they were as can be seen in figure 3.23.

The Average Spearman's Rank Correlation Coefficient was 0.14 and the Kendall's Tau was 0.096 with a *p-value* of 0.712.

### 3.4.7 Discussion

The first thing to address is that data quality is still an elusive concept since it has a contextual dimension and the quality of the record depends on the usage of the information. For example, data aimed at primary usage and day-to-day healthcare decisions about a patient will have different requirements regarding the importance of some variable or completeness of information very different from data needed to create summary statistics for key performance indicators extraction. Moreover, the data is still very vendor-specific. Even though we used an interoperability standard, the semantic layer, more connected with terminology is still lacking. This is an issue to be addressed in order to improve the interoperability of the standard. Moreover, we do not know how the training done with this data is generalizable to other vendors. One opportunity arises of mapping all of this data to a widely used terminology like SNOMED CT or LOINC. Nevertheless, the usage of FHIR and the fact that the data is mapped to a standard terminology, makes it easier

Figure 3.12: Model score for newly seen data

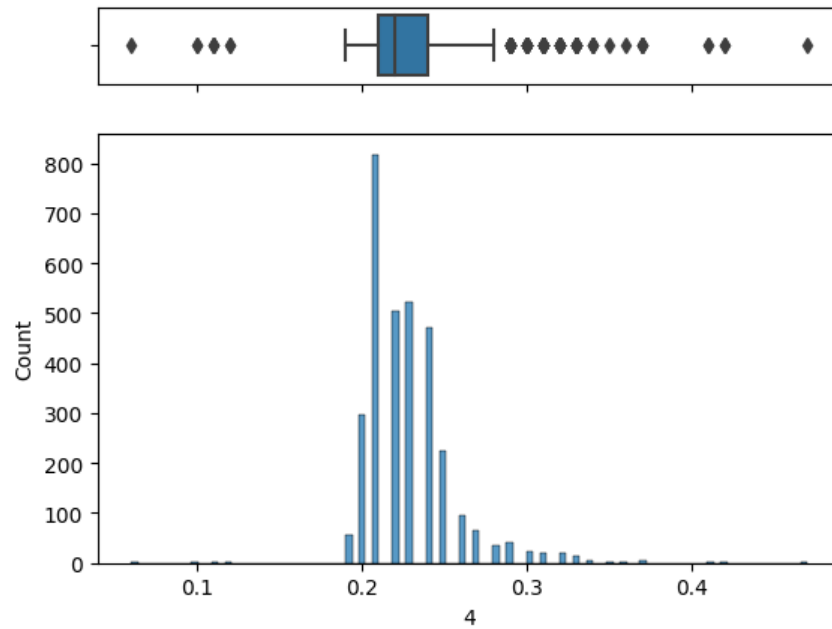
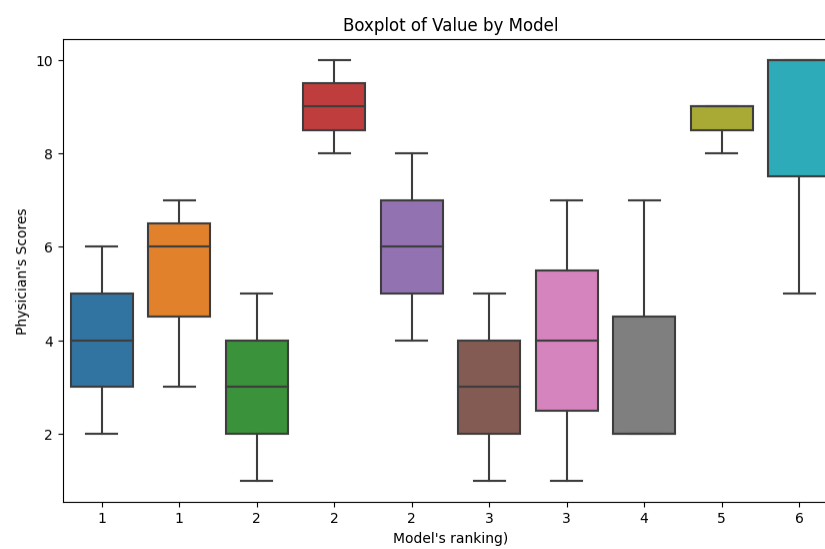


Figure 3.13: Comparasion of clinical assessment of records with the model



to use the data in other systems and to compare the results with other studies. Furthermore, being available freely and online makes it easier to understand how to map vendor-specific datasets to the model and use it in other contexts. Regarding the model, the usage of explainable methodologies like outlier-tree and transparent models like Bayesian networks are vital for clinical application. Since we use a single model to classify possible errors in the records, the ability to try to show clinicians why that value was tagged is of uttermost importance in order to get feedback and action from humans. From the experience gathered with the study, we believe that a weaker but transparent model could have more impact than better performant but opaque ones. If explainability and interpretability are important for any ML problem, this need only increases when we are dealing with such subjective concepts as data quality.

Regarding the clinical evaluation, we found out that asking clinicians to purely assess the quality of a record in an ehr is not an easy task. We found out that for a proper assessment, a context and objective must be defined in order to make the evaluation more objective. Moreover, the ranking methodology, even though is very useful for comparing with the model, it is not easy for clinicians to order 10 records when some of them have the "same level" of quality. This is a very important aspect to take into account when designing the evaluation of data quality. Probably a categorical evaluation of yes/no could be more useful and then compare with the ordering of the model and define thresholds based on that. These reasons are probably the cause behind the great variability between clinicians and between clinicians and the model. However, we do see a tendency for a higher agreement of the worst quality results than the best ones. This result suggests that the system may not be suitable for ranking good-quality records and clinicians are also not able. However, it could be useful to alert for low-quality ones, which is also a very important task with a great impact on the quality of the data. These findings are supported not only by 3.23 but also by figure 3.12 where we can add some threshold for the need for a human review. From the preliminary data in the questionnaires and looking at the graph, we believe that a threshold of around 0.3 could be a good starting point. However, this is a very subjective decision, and it should take into account the context and the objective of the evaluation. For example, if the objective is to use the data for research, a higher threshold could be used. On the other hand, if the objective is to use the data for day-to-day clinical decisions, a lower threshold could be used.

For the next steps, we believe a research path could be of identifying contexts for applying data quality checks like primary usage, research purposes, and aggregated analysis for decision-making among others could help better target those contexts and the importance of each variable for those use cases. This could be interesting to add to the tool in order to weigh the different variables according to the context.

### **3.4.8 Conclusion**

This work is still an early draft of a production-ready tool. However, we feel the work done is already a valuable insight into how to use data quality frameworks and several statistical tools in order to assess ehr data quality. This is a fundamental process not only to guarantee the quality of data for primary usage on a day-to-day but also for securing quality for secondary analysis and

usage. We believe the fact that we created an interoperable tool that was trained on real obstetrics data from 9 different hospitals and has the ability to provide a single score for a clinical record can help institutions, academics, and ehr vendors implement data quality assessment tools in their own systems and institutions.

For the next steps, we would like to further evaluate the score and its relationship with clinical usefulness. This would also include a further assessment of a threshold for the score for defining a record that would require human attention.

### **3.5 Leveraging Distributed systems in healthcare: is it advisable?**

This section is based on the paper entitled "Evaluating distributed-learning algorithms on real-world healthcare data". This paper was focused on the fact that access to healthcare data is often labourious and time-consuming. So we evaluated the distributed paradigm to its gold-standard, the centralized paradigm. We used 9 real-world datasets of obstetrics EHRs and compared the performance of several ML algorithms in both paradigms. We concluded that the distributed paradigm is a valid alternative to the centralized paradigm, with the added benefit of not requiring heavy data sharing.

#### **3.5.1 Introduction**

As the use of AI is increasing in the healthcare space [161], increased demand for ethical usage of personal patient data is occurring as well [39]. This has been happening both on the governmental side, with several regulations passed to protect citizens' data and personal information (such as GDPR in the EU [9] and HIPAA in the USA [138]), and on the public side, with an increased concern with continuous data breaches across institutions [5]. So, we are now faced with a dilemma on a compromise between what is possible to do with the available data and what should be done regarding patient privacy [207]. This is the main reason why health institutions implement burdensome processes and methodologies for sharing patient data, often costing a great deal of time, money, and human resources, seldomly overtaking the ideal time frame for analysing such data. Due to these privacy concerns, the traditional method for using data in healthcare is, nowadays, by focusing on data from a single institution in order to predict or infer something regarding those patients; this could be understood as local learning. This approach has some drawbacks, namely data quantity, data quality and possible class imbalance [157], never quite raising into its full potential for promoting best healthcare practices [215, 221, 205] with data sharing between institutions. In order to overcome this issue, there are a few, more complex, systems that aggregate data from several institutions, so more robust algorithms could be trained. However, this globally centralised aggregation of data encompasses a very important data breach hazard.

This is the setting where distributed learning could create a greater impact. A halfway point between local and centralised learning is where we train several models, one in each institution (or silo), and where the sole information that leaves the premises is a trained model or its metadata. A distributed model is built as the aggregation of all the local models, consequently aiming to

create a model similar to one globally trained with all the data in a centralised server. However, the distributed model never contacted with any data, only the local models did. This provides the opportunity to create better models, improve data protection, reduce training time and cost and provide better scaling capabilities [96].

There are already some implementations of distributed systems in the healthcare space, but we lack a robust understanding of how these models behave with real data, when compared with the classical models built with all the aggregated data. Additionally, the main issues regarding the development and implementation of such systems in healthcare are still elusive. So we aim to understand how distributed mechanisms behave compared to using all data in the healthcare space and if they are a suitable replacement for traditional machine-learning pipelines. The contributions of this paper are:

- Understand how to address the lack of data quality of real-world data regarding distributed model creation;
- Evaluate a distributed model against its local counterparts;
- Measure the prediction performance difference between a distributed model and a centralised one;
- Open a research path for using distributed models to predict several target variables in obstetrics clinical research.

### 3.5.2 Theoretical background and Related Work

Distributed learning [192] can be understood as training several models in a different setting and then aggregating them as a whole. There are two main branches of these approaches, distinguishable by the existence of a central orchestrator server: federated learning where such an entity exists, and peer-to-peer (or swarm) [207] learning where it does not. Even though distributed learning has been receiving a lot of attention recently, only some of its concepts have been focused on, mainly distributed-deep learning with a federated learning approach [216, 112]. These methods use the strength of neural networks and several algorithms like federated averaging to create distributed models capable of handling complex data like text, sound, or image [151]. However, considering that there are great amounts of information, especially in healthcare, stored as tabular data [11, 63, 144] and that neural networks are often not the best tool for such data structures [30], there is a lack of knowledge in the traditional machine learning techniques in a distributed manner.

Nevertheless, there have been some health-related distributed machine-learning projects successfully implemented, such as euroCAT [61] which implemented an infrastructure across five clinics in three countries. SVM models were used to learn from the data distributed across the five clinics. Each clinic has a connector to the outside where only the model's parameters are passed to the central server which acts as a master deployer regarding the model training with the radiation oncology data. Also, ukCAT [152] did similar work, with an added centralised database in the middle, but the training being done with a decentralized system.

Finally, a few works have explored the evaluation of models in a distributed manner, for example, comparing centralised machine learning, distributed machine learning and federated learning

on MNIST dataset [37]. Also, works that evaluate federated learning on MNIST, MIMIC-III and PhysioNet ECG datasets, but not in comparison with other methods [111]. The work by Tuladhar and colleagues [192] uses healthcare images and/or public and curated datasets. As far as we know, this is the first time a distributed machine learning evaluation is done with real-world clinical data from several different data sources.

### 3.5.3 Materials

Clinical data was gathered from nine different Portuguese hospitals regarding obstetric information, pertaining to admissions from 2019 to 2020. This originated nine different files representing different sets of patients but with the same features associated to them. The software for collecting data was the same in every institution (although different versions existed across hospitals) - ObsCare. The data columns are the same in every hospital's database. Each hospital was considered a silo and summary statistics of the different silos are reported in the tables 3.6 and 3.7. The data dictionary is in appendix A.1.

Table 3.6: Silos overview. categorical columns have a snippet of the most used category and a percentage. Continuous variables have a mean and standard deviation. Abbreviation meaning in the appendix. The last row is the number of patients. \* columns were used as target.

Column	Silo 1	Silo 2	Silo 3	Silo 4	Silo 5	Aggr.
IA*	31.1 <b>5.7</b>	30.7 <b>5.6</b>	31.1 <b>5.9</b>	31.1 <b>6.3</b>	31.3 <b>5.6</b>	31.1 <b>5.6</b>
GS*	a,rh.. <b>40%</b>	a,rh.. <b>40%</b>	a,rh.. <b>39%</b>	o,rh.. <b>38%</b>	a,rh.. <b>41%</b>	a,rh.. <b>40%</b>
PI	66.4 <b>14.4</b>	66.1 <b>13.5</b>	65.5 <b>14.1</b>	65.5 <b>14.1</b>	65.5 <b>14.4</b>	66.0 <b>14.1</b>
PAI	81.4 <b>14.9</b>	79.5 <b>14.5</b>	78.0 <b>15.2</b>	79.6 <b>16.3</b>	78.3 <b>14.2</b>	78.8 <b>14.5</b>
IMC*	25.2 <b>8.6</b>	25.2 <b>6.2</b>	25.0 <b>5.3</b>	25.0 <b>8.9</b>	24.9 <b>7.8</b>	25.1 <b>7.0</b>
CIG	Null <b>84%</b>	Null <b>85%</b>	Null <b>87%</b>	Null <b>90%</b>	Null <b>88%</b>	Null <b>88%</b>
APARA	Null <b>45%</b>	Null <b>41%</b>	1.0 <b>37%</b>	Null <b>42%</b>	1.0 <b>35%</b>	Null <b>39%</b>
AGESTA*	1 <b>41%</b>	1.0 <b>43%</b>	1.0 <b>39%</b>	1 <b>39%</b>	1 <b>43%</b>	1.0 <b>42%</b>
EA*	Null <b>75%</b>	Null <b>60%</b>	Null <b>75%</b>	Null <b>67%</b>	Null <b>45%</b>	Null <b>60%</b>
VA	Null <b>90%</b>	Null <b>80%</b>	Null <b>89%</b>	Null <b>93%</b>	Null <b>55%</b>	Null <b>77%</b>
FA	Null <b>99%</b>	Null <b>83%</b>	Null <b>94%</b>	Null <b>96%</b>	Null <b>60%</b>	Null <b>83%</b>
CA*	Null <b>88%</b>	Null <b>73%</b>	Null <b>86%</b>	Null <b>90%</b>	Null <b>62%</b>	Null <b>75%</b>
TG*	espo.. <b>62%</b>	espo.. <b>90%</b>	espo.. <b>85%</b>	espo.. <b>63%</b>	espo.. <b>89%</b>	espo.. <b>85%</b>
V	s <b>99%</b>	s <b>92%</b>	s <b>99%</b>	s <b>94%</b>	s <b>99%</b>	s <b>98%</b>
NRCPN*	7.3 <b>4.7</b>	7.0 <b>6.4</b>	6.4 <b>3.9</b>	5.5 <b>3.6</b>	10.5 <b>5.1</b>	8.4 <b>5.1</b>
VP	Null <b>82%</b>	Null <b>85%</b>	Null <b>81%</b>	Null <b>79%</b>	Null <b>73%</b>	Null <b>76%</b>
VCS	s <b>61%</b>	s <b>53%</b>	s <b>78%</b>	s <b>50%</b>	s <b>70%</b>	s <b>68%</b>
VNH	s <b>88%</b>	s <b>76%</b>	s <b>81%</b>	Null <b>52%</b>	s <b>71%</b>	s <b>69%</b>
B	Null <b>95%</b>	Null <b>78%</b>	Null <b>90%</b>	Null <b>97%</b>	Null <b>81%</b>	Null <b>83%</b>
AA	Null <b>89%</b>	Null <b>78%</b>	apr... <b>52%</b>	Null <b>96%</b>	Null <b>71%</b>	Null <b>73%</b>
BS	Null <b>98%</b>	Null <b>79%</b>	Null <b>97%</b>	Null <b>86%</b>	Null <b>97%</b>	Null <b>95%</b>
BC	Null <b>99%</b>	Null <b>83%</b>	Null <b>99%</b>	Null <b>87%</b>	Null <b>97%</b>	Null <b>97%</b>
BDE	Null <b>99%</b>	Null <b>83%</b>	Null <b>99%</b>	Null <b>88%</b>	Null <b>97%</b>	Null <b>97%</b>
BDI	Null <b>99%</b>	Null <b>83%</b>	Null <b>99%</b>	Null <b>87%</b>	Null <b>97%</b>	Null <b>96%</b>
BE	Null <b>99%</b>	Null <b>83%</b>	Null <b>99%</b>	Null <b>87%</b>	Null <b>97%</b>	Null <b>96%</b>
BP	Null <b>99%</b>	Null <b>83%</b>	Null <b>99%</b>	Null <b>87%</b>	Null <b>98%</b>	Null <b>97%</b>
IGA*	38.1 <b>3.5</b>	38.8 <b>2.2</b>	38.9 <b>1.6</b>	38.8 <b>2.4</b>	38.6 <b>2.1</b>	38.7 <b>2.2</b>
TPEE	Null <b>70%</b>	Null <b>75%</b>	Null <b>65%</b>	Null <b>64%</b>	Null <b>60%</b>	Null <b>65%</b>
TPEI	Null <b>98%</b>	Null <b>84%</b>	Null <b>93%</b>	Null <b>92%</b>	Null <b>99%</b>	Null <b>93%</b>
RPM	Null <b>91%</b>	Null <b>94%</b>	Null <b>89%</b>	Null <b>92%</b>	Null <b>85%</b>	Null <b>88%</b>
DG*	Null <b>88%</b>	Null <b>90%</b>	Null <b>90%</b>	Null <b>91%</b>	Null <b>90%</b>	Null <b>89%</b>
TP*	part.. <b>43%</b>	part.. <b>53%</b>	part.. <b>44%</b>	part.. <b>52%</b>	part.. <b>49%</b>	part.. <b>51%</b>
ANP	cefá.. <b>92%</b>	cefá.. <b>94%</b>	cefá.. <b>95%</b>	cefá.. <b>95%</b>	cefá.. <b>94%</b>	cefá.. <b>94%</b>
TPNP*	espo.. <b>53%</b>	espo.. <b>52%</b>	espo.. <b>58%</b>	espo.. <b>62%</b>	espo.. <b>62%</b>	espo.. <b>53%</b>
SGP*	38.5 <b>2.8</b>	38.9 <b>2.0</b>	39.1 <b>1.7</b>	39.0 <b>2.3</b>	38.9 <b>2.0</b>	38.9 <b>2.0</b>
GR*	1 <b>22%</b>	1 <b>20%</b>	1 <b>24%</b>	Null <b>81%</b>	1 <b>28%</b>	1 <b>24%</b>
N (total)	8039	8566	4989	2364	18177	80874

### 3.5.4 Methods

Data was prepossessed with the removal of features with high missing rates (> 90% in all silos). All missing value representations were standardized. The imputation process was done using the



Table 3.7: Silos overview part 2. categorical columns have a snippet of the most used category and a percentage. Continuous variables have a mean and standard deviation. Abbreviation meaning in the appendix. The last row is the number of patients. \* columns were used as target.

Column	Silo 6	Silo 7	Silo 8	Silo 9	Aggr.
IA	31.3 <b>5.2</b>	31.4 <b>5.4</b>	31.5 <b>5.6</b>	30.1 <b>5.6</b>	31.1 <b>5.6</b>
GS	a,rh.. <b>42%</b>	a,rh.. <b>39%</b>	a,rh.. <b>40%</b>	a,rh.. <b>42%</b>	a,rh.. <b>40%</b>
PI	65.6 <b>13.5</b>	66.0 <b>13.7</b>	65.6 <b>14.1</b>	67.4 <b>14.6</b>	66.0 <b>14.1</b>
PAI	77.7 <b>13.4</b>	79.2 <b>14.7</b>	76.7 <b>13.0</b>	83.1 <b>15.2</b>	78.8 <b>14.5</b>
IMC	24.9 <b>5.1</b>	24.9 <b>7.0</b>	24.8 <b>8.0</b>	25.7 <b>5.6</b>	25.1 <b>7.0</b>
CIG	Null <b>91%</b>	Null <b>91%</b>	Null <b>86%</b>	Null <b>90%</b>	Null <b>88%</b>
APARA	1.0 <b>38%</b>	Null <b>43%</b>	Null <b>41%</b>	Null <b>43%</b>	Null <b>39%</b>
AGESTA	1.0 <b>44%</b>	1 <b>43%</b>	1.0 <b>42%</b>	1.0 <b>40%</b>	1.0 <b>42%</b>
EA	Null <b>59%</b>	Null <b>61%</b>	Null <b>69%</b>	Null <b>61%</b>	Null <b>60%</b>
VA	Null <b>79%</b>	Null <b>82%</b>	Null <b>88%</b>	Null <b>82%</b>	Null <b>77%</b>
FA	Null <b>82%</b>	Null <b>86%</b>	Null <b>94%</b>	Null <b>89%</b>	Null <b>83%</b>
CA	Null <b>69%</b>	Null <b>75%</b>	Null <b>85%</b>	Null <b>78%</b>	Null <b>75%</b>
TG	espo.. <b>88%</b>	espo.. <b>85%</b>	espo.. <b>86%</b>	espo.. <b>93%</b>	espo.. <b>85%</b>
V	s <b>97%</b>	s <b>99%</b>	s <b>98%</b>	s <b>99%</b>	s <b>98%</b>
NRCPN	6.8 <b>4.0</b>	7.7 <b>3.2</b>	9.3 <b>4.5</b>	8.9 <b>5.5</b>	8.4 <b>5.1</b>
VP	Null <b>68%</b>	Null <b>74%</b>	Null <b>71%</b>	Null <b>78%</b>	Null <b>76%</b>
VCS	Null <b>53%</b>	s <b>87%</b>	s <b>63%</b>	s <b>87%</b>	s <b>68%</b>
VNH	Null <b>62%</b>	s <b>63%</b>	s <b>69%</b>	s <b>83%</b>	s <b>69%</b>
B	Null <b>90%</b>	Null <b>53%</b>	Null <b>93%</b>	Null <b>82%</b>	Null <b>83%</b>
AA	Null <b>84%</b>	apr... <b>61%</b>	Null <b>89%</b>	Null <b>74%</b>	Null <b>73%</b>
BS	Null <b>99%</b>	Null <b>98%</b>	Null <b>99%</b>	Null <b>95%</b>	Null <b>95%</b>
BC	Null <b>100%</b>	Null <b>100%</b>	Null <b>100%</b>	Null <b>97%</b>	Null <b>97%</b>
BDE	Null <b>100%</b>	Null <b>100%</b>	Null <b>100%</b>	Null <b>97%</b>	Null <b>97%</b>
BDI	Null <b>100%</b>	Null <b>100%</b>	Null <b>99%</b>	Null <b>97%</b>	Null <b>96%</b>
BE	Null <b>100%</b>	Null <b>100%</b>	Null <b>100%</b>	Null <b>97%</b>	Null <b>96%</b>
BP	Null <b>100%</b>	Null <b>100%</b>	Null <b>100%</b>	Null <b>97%</b>	Null <b>97%</b>
IGA	38.7 <b>1.8</b>	39.0 <b>2.0</b>	38.6 <b>2.1</b>	38.8 <b>1.9</b>	38.7 <b>2.2</b>
TPEE	Null <b>65%</b>	Null <b>64%</b>	Null <b>65%</b>	Null <b>63%</b>	Null <b>65%</b>
TPEI	Null <b>92%</b>	Null <b>86%</b>	Null <b>87%</b>	Null <b>94%</b>	Null <b>93%</b>
RPM	Null <b>85%</b>	Null <b>84%</b>	Null <b>90%</b>	Null <b>94%</b>	Null <b>88%</b>
DG	Null <b>92%</b>	Null <b>88%</b>	Null <b>90%</b>	Null <b>87%</b>	Null <b>89%</b>
TP	part.. <b>54%</b>	part.. <b>52%</b>	part.. <b>48%</b>	part.. <b>59%</b>	part.. <b>51%</b>
ANP	cefá.. <b>93%</b>	cefá.. <b>94%</b>	cefá.. <b>95%</b>	cefá.. <b>94%</b>	cefá.. <b>94%</b>
TPNP	espo.. <b>64%</b>	Null <b>100%</b>	espo.. <b>50%</b>	espo.. <b>65%</b>	espo.. <b>53%</b>
SGP	38.8 <b>1.8</b>	39.2 <b>1.7</b>	38.7 <b>2.0</b>	39.0 <b>1.6</b>	38.9 <b>2.0</b>
GR	1 <b>27%</b>	1 <b>25%</b>	1 <b>21%</b>	3 <b>27%</b>	1 <b>24%</b>
N (total)	12002	8258	6693	11786	80874

mean value (for continuous variables) or a new category (NULLIMP) for categorical variables. All categories were encoded as numbers using a previous mapping created based on all possible categories in all silos. Even though an ordinal relationship is created among features, we believe that

since we are applying this methodology to all datasets, which will be the source for all tests (local, distributed and centralised), that fact may be ignored. When training classification models, all of the target variable classes must be known at that moment and should be present in each split of the cross-validation. So, when assessing the training dataset, low-frequency target classes ( $n < 25$ ) were up-sampled with Synthetic Minority Oversampling Technique (SMOTE) [40] and missing target classes were addressed with dummy rows creation by the imputation of the mean for continuous variables and mode for categorical variables (per silo). These preprocessing mechanisms were applied in each run and for each target. The distributed model was an ensemble of models from each silo on a weighted soft-voting basis, defining weights and thresholds based on the training set scores. All procedures were coded in Python 3.9.7 with the usage of the scikit-learn library [146] and mlxtend library [159]. This study received Institutional Review Board approval from all hospitals included in this study with the following references: CHUSJ; 08/2021, CHBV; 12-03-2021, ULSM; 39/CES/JAS, HSOG; 85/2020, CHTS; 43/2020, CHVNGE; 192/2020, CHEDV; CA-371/2020-0t\_MP/CC, ULSAM; 11/2021.

### 3.5.4.1 Model Performance Evaluation

Local models were built with each silo's data. The distributed model was built as an ensemble of all the local models with weighted averaging. The centralised model was trained with a training dataset from all the silos combined. All models were built for a certain outcome variable with cross-validation and then compared, over 10 stochastic runs, with evaluation being performed on a test set held out from each silo. The metrics used for classification models were Weighted Area Under the Receiver Operating Characteristic Curve (AUROC) computed as One-versus-Rest, Weighted AUPRC. The metrics for regression models were Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE). The algorithm is shown in the algorithm 2. This rendered over 1000 different combinations.

### 3.5.4.2 Model Training

To avoid pitfalls of inductive bias from a certain learning strategy, we learned six different models (i) Decision Trees, (ii) Bayesian methods, (iii) a logistic regression model with Stochastic Gradient Descent, (iv) KNN, (v) AdaBoost and (vi) Multi-layer Perceptron. The decision was to create diversity in the models used, in order to assess if the training methodology could have an impact on distributed model creation. Nineteen features were used as target outcomes. These features were selected by filtering by the percentage of null (below 50%). For categorical outcomes, thirteen were selected. For continuous variables, six were selected. Details can be seen in tables 3.6 and 3.7.

After all the data was collected, we used the standard independent 2-sample T-test to check if the differences were significant with a  $\alpha$  of 0.05. We did the comparison between the distributed model and sequentially the centralised and correspondent local model across all algorithms.

Pre-process all silos (null standardization, imputation, encoding);

**for** *target in target list* **do**

    Create a centralised model with all the data with a 2x10 Cross-Validation

    Create distributed (ensemble of all models) model with:

**for** *silo in imputed silos* **do**

- Train-Test Split (80:20)
- check for low frequency or nonexistent labels in train set
- train local model with hyper-parameter tuning with 2x10 CV
- define weights based on scores in the train set

**end**

**for** *n in 10 repetitions* **do**

**for** *silo in imputed silos* **do**

- Train-Test Split (80:20)
- train local model with hyper-parameter tuning with 2x10 CV
- predict local on the test set
- predict distributed on the test set
- predict centralised on the test set

**end**

**end**

**end**

**Algorithm 2:** Creation and evaluation of the 3 different models

### 3.5.5 Results

Table 3.8 shows the aggregated metrics for AUROC, AUPRC, RMSE and MAE for distributed, centralised and local models predicting capabilities on each silo. The data refers to the mean of the metric values for all columns tested as targets for all methods and all silos. We also calculated the 95% confidence interval for each model (local and distributed per silo in order to assess how well the distributed model would work as opposed to the local one per silo. We also calculated the p Value for the means.

Table 3.8: Comparison for the centralised model, distributed model and local model (Mean for all model and all columns). Bold for  $P$  value below 0.05.

		M	SD	95% CI	$P$
AUPRC	distributed	0.691	0.216	(0.686, 0.696)	-
	centralised	0.706	0.225	(0.701, 0.711)	<b>1.10e-17</b>
	local	0.659	0.220	(0.654, 0.665)	<b>4.71e-05</b>
AUROC	distributed	0.723	0.182	(0.718, 0.727)	-
	centralised	0.729	0.180	(0.725, 0.734)	<b>2.98e-26</b>
	local	0.692	0.164	(0.688, 0.695)	<b>2.48e-02</b>
MAE	distributed	2.370	1.608	(2.315, 2.425)	-
	centralised	2.365	1.923	(2.298, 2.431)	<b>2.23e-04</b>
	local	2.527	1.799	(2.465, 2.589)	9.01e-01
RMSE	distributed	21.171	46.078	(19.584, 22.757)	-
	centralised	19.839	28.645	(18.853, 20.826)	<b>2.92e-02</b>
	local	23.771	49.776	(22.057, 25.485)	1.63e-01

### 3.5.6 Discussion

The imputation process was done using the mean value (for continuous variables) or a new category (NULLIMP) for categorical variables. All categories were encoded as numbers using a previous mapping created based on all possible categories in all silos. Even though an ordinal relationship is created among features, we believe that since we are applying this methodology to all datasets, which will be the source for all tests (local, distributed and centralised), that fact may be ignored. When training classification models, all of the target variable classes must be known at that moment and should be present in each split of the cross-validation. So, when assessing the training dataset, low-frequency target classes ( $n < 25$ ) were up-sampled with SMOTE [40] and missing target classes were addressed with dummy rows creation by the imputation of the mean for continuous variables and mode for categorical variables (per silo). These preprocessing mechanisms were applied in each run and for each target. The distributed model was an ensemble of models from each silo on a weighted soft-voting basis, defining weights and thresholds based on

Figure 3.14: Heatmap of classification algorithm and silo vs Target variable and model type. Value is the AUROC mean of all 10 experiments. Y axis is the algorithm and silo. X axis is Target variable and Method.

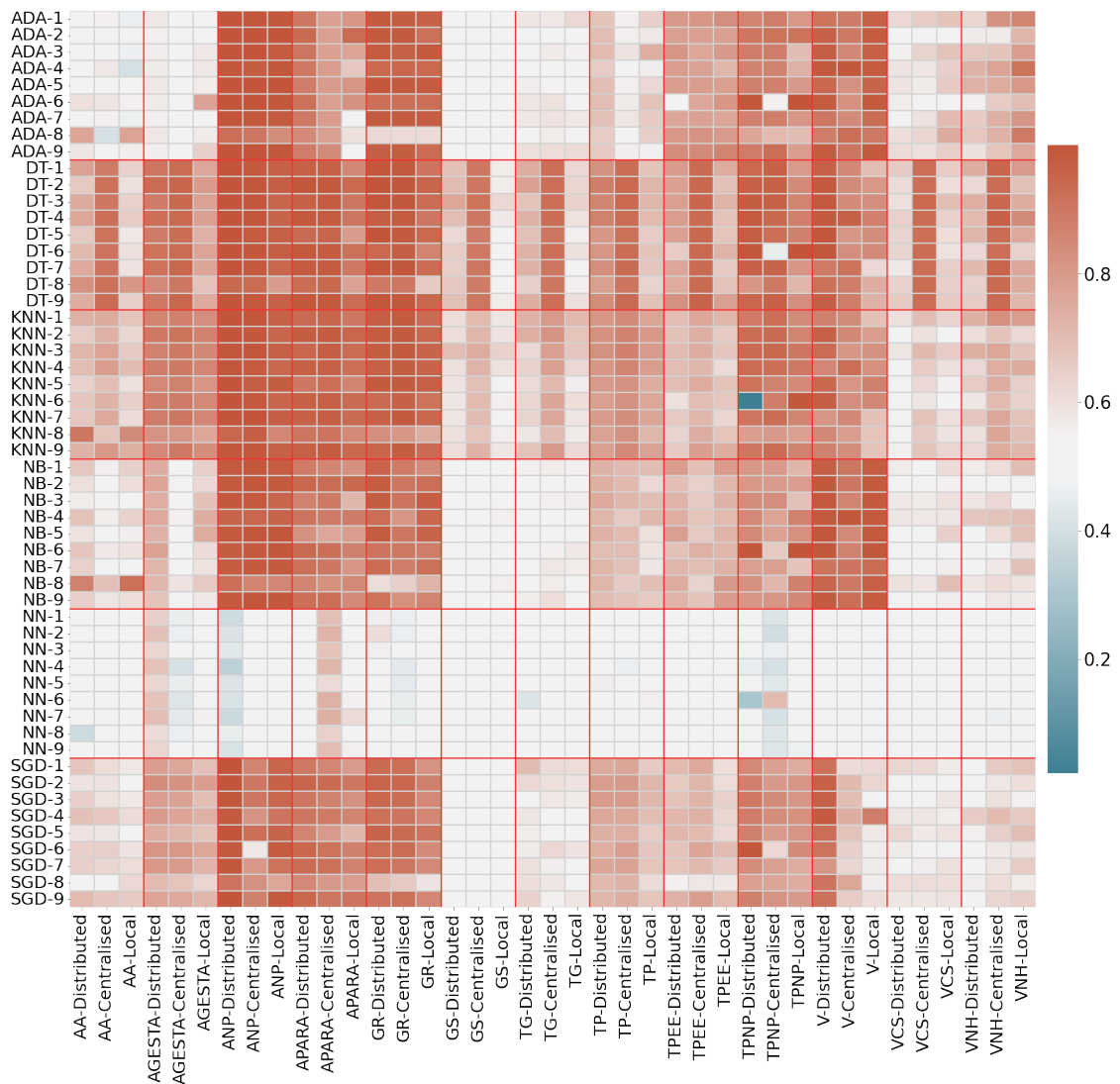


Figure 3.15: Heatmap of regression algorithm and silo vs Target variable and model type. Value is the MAE mean of all 10 experiments. The y axis is the algorithm and silo. X axis is Target variable and Method.



Table 3.9: Hypothesis testing of Distributed versus Centralised and local for every test. Each cell is the total of distributed model when compared with centralised model (row) and local model (column). (> for better, = for non significance and < for worse)

		> Local	= Local	< Local	Total
SGD	> Centralised	72	14	9	95
	= Centralised	14	17	6	37
	< Centralised	11	11	17	39
NN	> Centralised	44	44	7	95
	= Centralised	2	33	2	37
	< Centralised	0	17	22	39
KNN	> Centralised	16	0	1	17
	= Centralised	10	2	1	13
	< Centralised	72	28	41	141
ADA	> Centralised	64	25	22	111
	= Centralised	5	12	10	27
	< Centralised	10	6	17	33
NB	> Centralised	51	19	34	104
	= Centralised	5	19	12	36
	< Centralised	3	4	24	31
DT	> Centralised	27	0	1	28
	= Centralised	8	0	0	8
	< Centralised	97	12	26	135
Total		511	263	252	1026

the training set scores. The first thing that is noticeable is the high scores achieved in our analysis which show that all algorithms in all forms (local, distributed and centralised) have a good grasp on ranking data (negative on the bottom and positives on the top of a scale) for classification or predicting the value for regression. We notice that distributed models have performance similar to their centralised counterparts. ~59% of all of the distributed models had similar or better performance than the centralised models. This suggests that a distributed model can be used to reliably infer information and does not compromise prediction performance when compared with the gold standard (centralised) while increasing privacy for the data owners.

Overall, our results suggest that it is possible to implement a distributed model without significantly losing information. However, there are still issues to be addressed. This methodology presents hurdles regarding categorical class handling. Firstly, all classes should be known first-hand and should be given to each model even if that silo in particular has no cases of that class. Secondly, low-frequency classes are also an issue to be addressed, since training the model with cross-validation will raise problems because each split should have all classes present. Our approach relied on sample creation for low and non-existent target classes. However, this approach is adding information to the model that is not originally there. The way we chose for minimising this issue was by creating dummy variables with median and mode imputations based only on the information in the dataset. Nevertheless, non-existent classes are impossible to address without prior information. These class problems could be partially tackled in production by implementing data management and governance procedures, namely data dictionaries. Still on data preprocessing, we applied ordinal encoding to the variables which will create a natural hierarchy between variables. One solution for this is to create binary columns for each class in each column. This will remove the hierarchy between classes but increase variable numbers and training time considerably.

Moreover, like in most secondary usage of data, other issues are important to keep in mind, even in such a controlled environment as this one. Even though the software is the same in every hospital, the clinical service is the same and the underlying data models are the same, the version of the software is not the same across all hospitals. This difference alone can alter the way each column is populated, mainly through front-end changes or label modification, among other aspects. Additionally, each hospital has its own workflows in practice that can also alter the way data is collected; changing timings or steps in a certain workflow can dramatically change the data acquisition and the reality it represents.

Another issue to consider is the path adopted to build the distributed model. In this case, it was decided to develop an ensemble of models with voting. However, other methods could have been employed, like parameter averaging, that should be tested as well. In particular, the usage of more robust neural networks could be assessed as well. We chose not to test state-of-the-art neural networks since the data volume was low for that use case and several papers have already demonstrated that neural networks are not the most suitable tool for tabular data [81, 29]. We chose to add MLPerceptron as a baseline for comparison with the remaining algorithms. The results show us that the performance was below the other algorithms, but in this concrete case, the problem



may reside in the architecture chosen and hyperparameters used in the Cross-validation. Despite this, a precise and thorough demonstration of this use case would be important to consider such scenarios.

Furthermore, the algorithm underlying the distributed model is of importance as well for its performance versus the centralised model. Figures 3.14 and 3.15 and table 3.9 show us that decision trees and K-nearest neighbours implemented in a centralised manner are consistently better than the distributed counterpart. Even though this improvement may have a relationship to the target variable (i.e figure 3.15 for IA and IGA variables), it is still an important fact to take into account when implementing such architectures. The performance of the models is also interesting to catch differences in silos. See silo 6 for TPNP (figure 3.14) where silo 6 consistently behaves differently than the rest. As for implementation, such a mechanism could be implemented in at least two manners; with a central orchestrator or without. The first one would assume a central point that would make a request to each silo for a prediction and then create the final prediction with the weighted averaging of each one. The second one would not require any additional platform and each silo would communicate with each of the others and receive the prediction and would create the final with their own. This implementation step would of course take into account variables that we were out of scope such as the communication between silos. Regarding the prediction capability as a whole, we found that this data is suitable to apply ML models in order to predict several clinical outcomes, with very good results for several target variables.

### 3.5.7 Conclusion

With this paper, it was possible to assess how well-distributed models can perform with real data, when compared to local models (trained with data from each silo) and global centralised models (trained with all data). These results show that an ensemble of models is able to fully grasp the specificity of the data, with performance similar to that of a model built with all the data. Even though the nature of the target and the silos can impact the performance, and several issues should be considered during the implementation phase, we are now fairly confident that distributed learning is a step forward regarding data privacy without loss of prediction performance. Finally, taken into account that the scores for several target variables are AUROC/AUPRC above 80% and MAE below 1, we will explore this further in a different work. We hope to be able to develop distributed models for predicting clinical outcomes like delivery type or Robson group, that could turn out useful in real-world clinical practice.

## 3.6 Can Institutions share their performance metrics without hesitation of retaliation?

This section is based on the paper entitled "Benchmarking institutions' health outcomes with clustering methods". This paper was focused on the fact that many healthcare institutions harbor reservations about openly sharing production metrics. One predominant concern is the potential

for retaliatory actions, be it from regulatory bodies, competitors, or the public. In this paper, we propose the application of a clustering methodology that allows institutions to compare performance metrics without disclosing the actual values. The method is based on clustering, which involves grouping health institutions' outcomes into a known number of clusters, allowing institutions to position themselves in a range of clusters without sharing the true means of their target data. The proposed method uses the K-means and K-modes clustering algorithms and was tested on data from real Electronic health records and public datasets. This approach provides a valid benchmark of hospital metrics and performances while protecting the privacy of participating institutions.

### 3.6.1 Introduction

Health institutions play a critical role in providing essential healthcare services to communities and ensuring that they operate efficiently and effectively is crucial. Benchmarking is a process that allows hospitals to compare their performance against that of other institutions, which can help identify areas of strength and weakness [181]. By analyzing and evaluating performance metrics, such as patient outcomes, operational efficiency, and financial management, hospitals can identify best practices and make data-driven decisions to improve their overall performance. It can also help hospitals identify and implement innovative practices that can lead to better patient care and improved staff satisfaction [91].

However, despite the numerous benefits of benchmarking, some hospitals may be hesitant to participate due to concerns about revealing weaknesses or being perceived as inferior to their peers. The fear of being judged or penalized for poor performance can sometimes lead hospitals to avoid sharing data, making it difficult to accurately assess their performance and identify areas for improvement. Privacy issues and concerns turn this opportunity into an even less desirable path [91]. To address these concerns, benchmarking initiatives often ensure the confidentiality and anonymity of data to encourage participation and foster trust among participating institutions. However, this is usually not enough. In 2019, as stated in the work of Villanueva et. al., [198], 26% of data-sharing initiatives are based on the aggregation of data and 24% are based on sharing data in closed consortia. Only 15% were based on open or controlled access.

To address concerns around privacy and confidentiality, we propose a new method of benchmarking based on clustering. This method involves grouping health institutions' outcomes into a known number of clusters, providing health institutions with the capability of positioning themselves in a range of clusters, without ever sharing the true means of their target data.

This approach to benchmarking not only addresses concerns around privacy and confidentiality. It has the potential to encourage greater participation in benchmarking initiatives, as hospitals can be assured of the anonymity and confidentiality of their data. By creating a more secure and private environment for benchmarking, hospitals can feel more comfortable sharing their data and participating in initiatives that can ultimately improve patient care and operational efficiency.

In conclusion, benchmarking is a crucial tool for hospitals to improve their performance and provide better care for their patients. While concerns around privacy and confidentiality may exist,

the clustering approach to benchmarking provides a more accurate assessment of hospital performance while protecting the privacy of participating institutions. By embracing benchmarking initiatives and leveraging new approaches to benchmarking, hospitals can continuously improve their operations and ensure they provide the highest quality of care possible. In this paper we propose:

- study how to implement clustering mechanism for benchmark
- address preprocessing issues for the raw data
- highlight pain points to deployment in the real world.

### 3.6.2 Rationale and Related Work

This work was initially suggested as a follow-up to a previous work of Rodrigues et al., [164] where clustering is applied to streaming data sources. We then thought if a similar approach could be applied to healthcare in order to be able to compare data distributions without ever knowing their real values of them. Clustering in healthcare is often used to create clusters of patients, taking into account a given set of characteristics. This is used to find possible groups of phenotype and be able to characterise populations given the centroids [203, 24]. It is also used as a method of detecting regularities and patterns in multi-omics data that reveal different molecular subtypes [134, 158]. It can also be used to create unsupervised models for facilitating the annotation of data for supervised models [124].

K-means [118, 180, 121] is an unsupervised clustering algorithm used to group data points into K distinct clusters based on their similarity. It is widely used in machine learning, data mining, and image segmentation. The algorithm works by randomly initializing K centroids (or cluster centres) and assigning each data point to the nearest centroid. Then, the centroids are moved to the mean of the points assigned to each cluster. This process is repeated until convergence, where the clusters no longer change.

The objective of K-means is to minimize the sum of squared distances between each data point and its assigned centroid, which is also called the within-cluster sum of squares (WCSS). The algorithm attempts to find the best K clusters that minimize the WCSS. However, choosing the right value of K can be challenging, and the algorithm may converge to a suboptimal solution. Therefore, K-means is often run multiple times with different initializations to find the best clustering solution. Despite its simplicity, K-means can be computationally expensive when dealing with large datasets, and it may not work well with non-linearly separable data or when the clusters have different shapes and sizes.

K-modes is another clustering algorithm similar to K-means, but it is designed to work with categorical data. Unlike K-means, which computes the mean of continuous variables, K-modes computes the mode (or the most frequent value) of categorical variables within each cluster. The algorithm works by randomly initializing K centroids and assigning each data point to the nearest centroid based on the number of matching categories. Then, the centroids are moved to the mode of the categories within each cluster. This process is repeated until convergence, where the clusters no longer change.

The objective of K-modes is to minimize the dissimilarity between the data points within each cluster, which is often measured by the Hamming distance, Jaccard distance, or other similarity measures. Like K-means, choosing the right value of K is critical, and the algorithm may converge to a suboptimal solution. Therefore, K-modes is often run multiple times with different initializations to find the best clustering solution. K-modes is particularly useful when dealing with data that have a large number of categorical variables or when the data contain missing values. However, like K-means, K-modes may not work well with non-linearly separable data or when the clusters have different shapes and sizes.

However, as far as we know, this is the first time clustering is tested for exchanging information privately.

### 3.6.3 Materials & Methods

#### 3.6.3.1 Method Overview

We used Python 3.9 to implement the mock example of such an use-case. The clustering was done with scikit-learn library [146]. The algorithm proposed is shown in algorithm 3.

```

for variable in silo do
  | initialize centroids;
end
while No convergence do
  |
  | • Send centroids to other silos
  |
  | • Receive other silo's information and add own centroids
  |
  | • Calculate new centroids
  |
  | • calculate score
end

```

**Algorithm 3:** Benchmarking with clustering

The method for assessing convergence is based on clustering metrics: the Rand Index (RI). This metric computes a similarity measure between two clusters by considering all pairs of samples and counting pairs that are assigned in the same or different clusters in the predicted and true clusters [90]. The raw RI score is:  $RI = (\text{number of agreeing pairs}) / (\text{number of pairs})$ . Furthermore, convergence must be obtained through several iterations to make sure it's stable, so a buffer period is also important. For the results section, we set the threshold as 0.9 and repetitions at 20.

In this paper, we propose to show how such an implementation could be done while addressing issues with data formats, types and preprocessing. So, we want to check if the encoding of categorical data affects the model and which method is better for encoding such variables. Additionally, we will try to understand if it is possible to create mechanisms for mixed data if categorical and continuous data must be used and evaluated separately and if so, through which mechanisms. We will test (1) continuous variables alone, and (2) encoded categorical variables as ordinal. We will

also test (3) K-modes and (4) K-means with the proportion of each category for categorical data. K-means was used as implemented in scikit-learn [146] and K-modes, as implemented by J. de Vos [60].

### 3.6.3.2 Data used

We used two types of data in this paper. One is simpler and available online from the UCI dataset library, namely, the heart disease dataset [95]. We made fairly simple preprocessing on that dataset, namely removing the "?" by filling with null and then imputing missing values by imputing the mean on continuous variables and mode on categorical ones. We then separated the data into 3 distinct silos at random to mimic different health institutions.

In order to use real data and address problems found in the wild, we used clinical data gathered from nine different Portuguese hospitals regarding obstetric information, pertaining to admissions from 2019 to 2020. This originated from nine different files representing different sets of patients but with the same features associated with them. The software for collecting data was the same in every institution (although different versions existed across hospitals) - ObsCare. The data columns are the same in every hospital's database. Each hospital was considered a silo for comparison.

### 3.6.4 Results

As for results, the data from heart disease rendered the figure 3.16. In this, we focused on continuous variables only. For easier reading, the data is as shown in the table 3.10. We used data from the real world to test if everything would work similarly, rendering the image 3.17. We added a binary category to show how meaningless the value turn in order to get any information out of it.

Figure 3.16: Clustering for 3 continuous variables with 3 silos and true centroids (S2) and true means (S2) for example purposes; The values were normalized for visualization purposes with MinMax

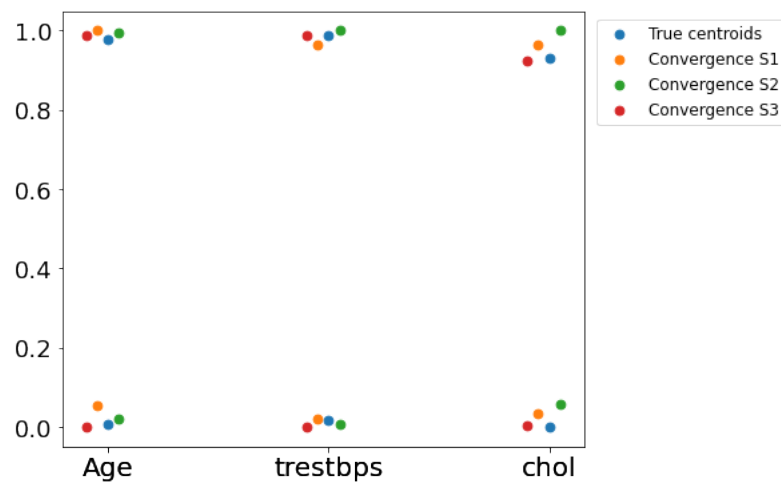
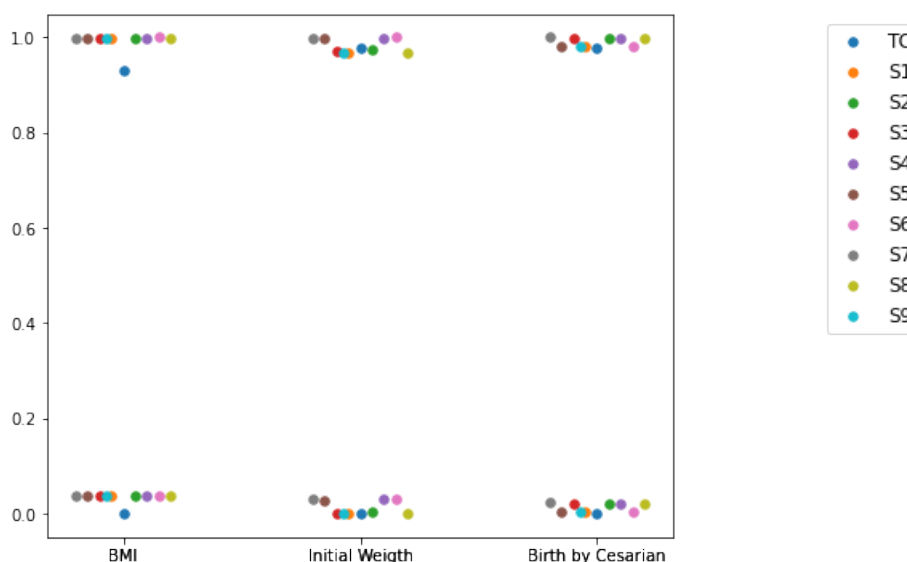


Table 3.10: Final Data points after convergence; S1, S2 and S3 are the centroids obtained in each silo (S) after convergence; True centroids are the centroids of the true means of all silos (TC)

	Age	trestbps	chol
S1	46.3 , 61.1	121.1 , 148.9	218.9 , 300.8
S2	45.8 , 61.0	120.7 , 149.9	220.9 , 304.0
S3	45.5 , 61.0	120.5 , 149.6	216.1 , 297.4
TC	45.6 , 60.8	121.0 , 149.6	215.8 , 297.9

Figure 3.17: Clustering for 3 variables with 9 silos and true centroids of the true means (TC); 2 continuous and 1 categorical one hot encoded, The values were normalised for visualisation purposes with MinMax



As before, the data is in table format in 3.11.

Then we experimented with categorical variables. Figure 3.18 shows the convergence of the silos with proportion data and K-means with that and with K-modes.

### 3.6.5 Discussion

As per the discussion, there are a few issues to be addressed. First as per data preprocessing. In order to cluster be obtained, the null data must be filled out. There are a few strategies to do so. One option is to eliminate records/rows with empty cells or impute data. Either is a possibility, with pros and cons but the capability of having a dataset where no null records are present across several features may be difficult to find in the wild, especially since there are often optional and conditional fields in most EHR. So imputation becomes more interesting, since it enables the usage of the whole dataset, even if biases are introduced. Mixed types of datasets are also an issue to be aware of. In this case, not only imputation but also encoding a categorical variable is a vital

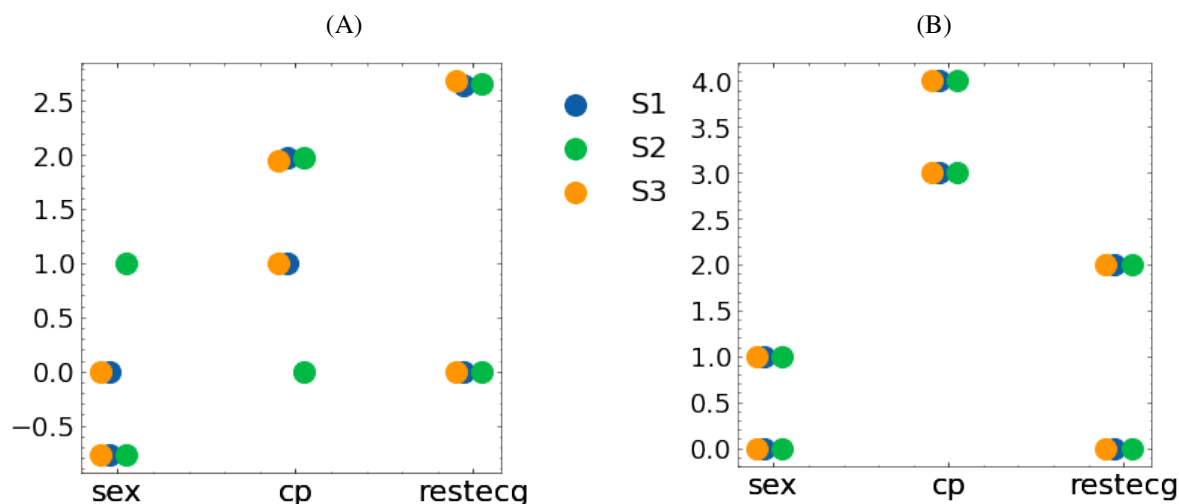
Table 3.11: Final Data points after convergence and true centroids of the true means of each silo (TC)

	BMI	Initial Weight	Birth by Cesarian
TC	24.9 , 383.1	60.5 , 85.0	0 , 1
S1	40.1 , 409.4	60.4 , 85.0	0.96 , -0.04
S2	40.1 , 410.4	61.7 , 86.3	0.99 , -0.01
S3	40.0 , 410.4	61.9 , 86.5	0.96 , -0.04
S4	40.6 , 411.3	61.9 , 86.5	0.96 , -0.04
S5	40.0 , 410.4	60.5 , 85.1	1.0 , 0.0
S6	40.1 , 409.3	60.4 , 84.9	1.0 , 0.0
S7	40.7 , 411.3	60.5 , 85.0	0.96 , -0.04
S8	40.0 , 410.4	86.5 , 61.9	1.0 , 0.0
S9	41.0 , 410.4	85.0 , 60.4	1.0 , 0.0

step to take in the preprocessing phase. There are usually two main methods of data encoding, ordinal encoding and binary encoding. The first one keeps a unique column as the original data but maps every category to an increasing natural number. This creates an ordering in the data, often a misrepresentation of reality, not only due to this hierarchy but only because it assumes the differences between ranks of the hierarchy are always the same (1). The second is related to expanding the number of columns into the number of categories and creating 0s and 1s for the category. In machine-learning terms, binary seems more suited to be applied, but for benchmarking purposes, both are below par in terms of interpretability. For categorical data, we found out that K-modes seem to fulfil the requirements in a better way, providing better interpretability and reasoning about the results. However, it should be noted that we applied K-modes in a multivariate fashion and K-means in a univariate fashion. Given that no percentage is provided, only the mode of the data, we believe it is still hard to get any real insight from the centroids. However, K-modes provides less information, since it only shows the top two categories. Which, for example. binary targets, provide little to no information. However, for larger categorical sets, the information provided could be better. Moreover, the number of centroids pretended could be more important as well. Agreeing on only 1 centroid would render the mode of the data provided by all silos, which could be more interesting. As for continuous data, the use of real data was insightful, since BMI had a few very big outliers around 300 and 400, which rendered centroids around that data. Even if not all silos had examples of these outliers, the ones that do have, pass that into the remaining. One possible workaround would be an addition of an extra cluster in order to catch possible outliers. However, this should be addressed in detail and assess how outliers could subvert the data from the silos and how to work around that.

As for the next steps, a few issues could be addressed in depth. Regarding imputation, it could

Figure 3.18: Clustering for 3 variables with 3 silos - (A) categorical variables with proportion with K-Means and (B) Categorical with K-modes



be interesting to understand how imputation, and which methods are more suitable to use for real-world scenarios. If the imputation of variables with a high null percentage influence significantly a centroid formation. Communication could be important as well. Which action is to be taken when a silo is "down" and does not send information to the remaining. Cluster information should be addressed as well. They need to be agreed upon beforehand in the scope of this paper. But if it could be selected by each silo? Would that be feasible or a convergence could be achieved? Finally, there is the question if there is the possibility of having leaks of true means across iterations by adversarial learning. At present time, we cannot be sure that the values are totally private, but then again, nothing is.

### 3.6.6 Conclusion

We believe that this work helps create the foundation for exchanging data across healthcare institutions without revealing the true data points. It could be useful for benchmarking and promoting a higher adoption rate. Even though there are still issues to be addressed, we think that the path is full of possibilities.

## 3.7 Leveraging data to assess treatment efficacy

This section is based on the paper entitled "Comparative Analysis of Palbociclib and Ribociclib: A real world data and Propensity Score-Adjusted Evaluation with endocrine therapy". This was a method of applying the knowledge of causality and transparent ML models in order to assess the real-world effect of two drugs for breast cancer. We started with traditional analysis and then moved to a more complex approach, using IPTW methods in order to further compare treatments.



### 3.7.1 Introduction

Currently, metastatic breast cancer is difficult to treat. Patients with Hormone Receptor-positive (HR+) and Human Epidermal Growth Factor Receptor 2-negative (HER2-) breast cancer, the most common subtype, typically undergo endocrine therapy. Therefore, new treatments can be very useful in improving quality of life, reducing toxicity, and decreasing scenarios of hormonal resistance. Medications from the group of cyclin-dependent kinase inhibitors appear as a potential improvement in the therapeutic approach to advanced breast cancer. Within this group, there are palbociclib, ribociclib and abemaciclib. Cyclin-dependent kinases 4 and 6 (CDK4/6) are responsible for regulating the cell cycle at the transition between the G1 and S phases. In many neoplasms, this cycle is deregulated, and it promotes uncontrolled cell proliferation. It is then possible for these medications to have better effectiveness. These medications were approved by INFARMED, I.P. after an analysis of the therapeutic value they offer. This decision was made based on data provided by clinical trials done with these medications. The MONALEESA [88, 176, 191] studies were used for ribociclib, PALOMA [196, 169, 71] for palbociclib, and MONARCH [75, 177] for abemaciclib. These studies focused on testing the hypothesis of treating CDK4/6 inhibitors in combination with an aromatase inhibitor or fulvestrant as an alternative to the gold standard. In these research findings, it was determined that there was a notable enhancement in effectiveness, supporting their application in clinical practice. However, this evaluation was based on clinical trials with very specific inclusion and exclusion criteria and in a highly controlled environment. It is then vital to study how these new molecules compare to current practice in terms of treatment effectiveness in a real-world setting. In the meticulously controlled setting of clinical trials, patient selection often skews towards relatively healthier individuals with fewer comorbidities. However, in real-world clinical practice, patients present a diverse range of health profiles, co-existing illnesses, and medication histories that may influence drug efficacy and safety. Real-world data, drawn from electronic health records, insurance claims databases, and patient registries, offers the advantage of reflecting a more heterogeneous patient population, thus potentially uncovering insights not readily apparent in clinical trial settings. Understanding the effectiveness and safety of CDK4/6 inhibitors in real-world conditions is crucial for tailoring more individualized treatment regimens, optimizing outcomes, and enhancing the quality of life for patients with HR+, HER2- breast cancer [84]. Nevertheless, observational studies have inherent limitations, such as confounding by indication, which can lead to biased estimates of treatment effects. To tackle this, there are causality-based assessments that can be employed in order to better estimate the causal effects of treatments. Incorporating statistical techniques like Inverse Probability of Treatment Weighting (IPTW) can play an essential role in enhancing the quality of real-world evidence by accounting for treatment selection bias and balancing observed covariates between treatment groups. IPTW, grounded in the framework of causal inference, allows for the mimicking of a randomized control trial-like setting within observational studies. By assigning weights to individual patients based on their propensity scores—the likelihood of receiving a particular treatment given a set of observed characteristics—analyses can achieve a balance between different treatment arms,

thereby reducing bias and confounding factors. Establishing causality, rather than mere association, is vital for the robust interpretation of real-world data. As we strive to understand the long-term impact, efficacy, and safety of CDK4/6 inhibitors in HR+, HER2- breast cancer, the rigorous application of IPTW and causal inference methods can substantially augment the validity of real-world findings, making them a more reliable basis for clinical decision-making [15, 17] So in this paper, we propose:

- To compare the effectiveness of the CDK4/6 inhibitors drug class in terms of progression-free survival (PFS) and overall survival (OS);
- Assess the Hazard Ratio of using the CDK4/6 inhibitors drug class in terms of PFS and OS.
- To compare the effectiveness of CDK4/6 inhibitors in combination with an aromatase inhibitor or fulvestrant with the current standard of care in terms of PFS and OS in patients with HR+, HER2- advanced breast cancer.
- assess the differences in effectiveness between the three CDK4/6 inhibitors in combination with an aromatase inhibitor or fulvestrant in terms of PFS and OS with causality principles in mind, especially the counterfactual theory and IPTW ;

### 3.7.2 Materials & Methods

#### 3.7.3 Study Design

This retrospective study was designed in 2022. The aim of the study was to evaluate the clinical benefit and long-term survival of patients with HR+/HER2- that started treatment with CDK4/6 inhibitors plus endocrine therapy in different lines of treatment between the 14th of March 2017 and the 31st of December 2021. The follow-up period was set until June 2022. Inclusion criteria: women and men, Hormone receptor-positive and HER2 negative in the primary tumor or metastatic site after biopsy. Exclusion criteria: Patients that had only ambulatory medication, and patients involved in clinical trials, diagnosed with other neoplasms or with active treatment during the study period. The comparison group was defined by a population of patients, that were treated with hormone therapy as first-line (due to bone metastases) between 2015 and 13 of March 2017.

The evaluation of effectiveness will involve overall survival and progression-free analysis. We will compare the three different cyclin-dependent kinase inhibitors in terms of efficacy in real-world patients and will also compare the effectiveness of this class of drug against traditional endocrine therapy.

#### 3.7.4 Data collection

All data were collected from original medical records from baseline to last visit or death. The data was collected from Instituto Português de Oncologia – Porto (IPO-P). Table 3.12 shows a comparison between the groups. Data included for population treated with CDK4/6 inhibitors plus endocrine therapy: demographic information, age at first diagnosis and age at the beginning of

Table 3.12: Descriptive statistics of cyclin-dependent kinase inhibitors group and endocrine therapy group. The Drug/combination refers to the actual drug or the combination for CDK4/6

	ET	Palbociclib	Ribociclib
	(N=43)	(N=246)	(N=106)
<b>Age at treatment start</b>			
Mean (SD)	60.1 (12.4)	59.2 (11.7)	58.2 (10.7)
Median [Min, Max]	62.0 [34.0, 85.0]	60.0 [28.0, 84.0]	58.0 [32.0, 79.0]
<b>Bone Only metastases</b>			
No	NA	161 (65 %)	74 (70 %)
Yes	NA	85 (35 %)	32 (30 %)
Missing	43 (100%)	0 (0%)	0 (0%)
<b>Visceral metastasis</b>			
No	NA	121 (49 %)	49 (46 %)
Yes	NA	125 (51 %)	57 (54 %)
Missing	43 (100%)	0 (0%)	0 (0%)
<b>Stage</b>			
I	3 (7 %)	22 (9 %)	7 (7 %)
II	20 (47 %)	75 (30 %)	22 (21 %)
III	11 (26 %)	74 (30 %)	18 (17 %)
IV	2 (5 %)	65 (26 %)	46 (43 %)
Missing	7 (16.3%)	10 (4.1%)	13 (12.3%)
<b>Drug/Combination</b>			
Anastrozol	3 (7 %)	NA	NA
Exemestane	4 (9 %)	NA	NA
Fulvestrant	5 (12 %)	180 (73 %)	10 (9 %)
Letrozol	31 (72 %)	66 (27 %)	96 (91 %)

treatment, clinical characteristics and performance status by Eastern Cooperative Oncology Group scale (ECOG), treatment line and treatment schema - CDK4/6 inhibitor and endocrine therapy, stage of cancer, site of metastases (bone, soft tissue, visceral, central nervous system-CNS with or without another site). Data included for the population treated with endocrine therapy as first-line: demographic information, age at first diagnosis and age at the beginning of treatment, clinical characteristics and performance status by Eastern Cooperative Oncology Group scale (ECOG), stage of the cancer.

For comparison purposes, we used palbociclib and ribociclib since we had a small number of patients treated with abemaciclib (12).

### 3.7.5 Statistical Analysis

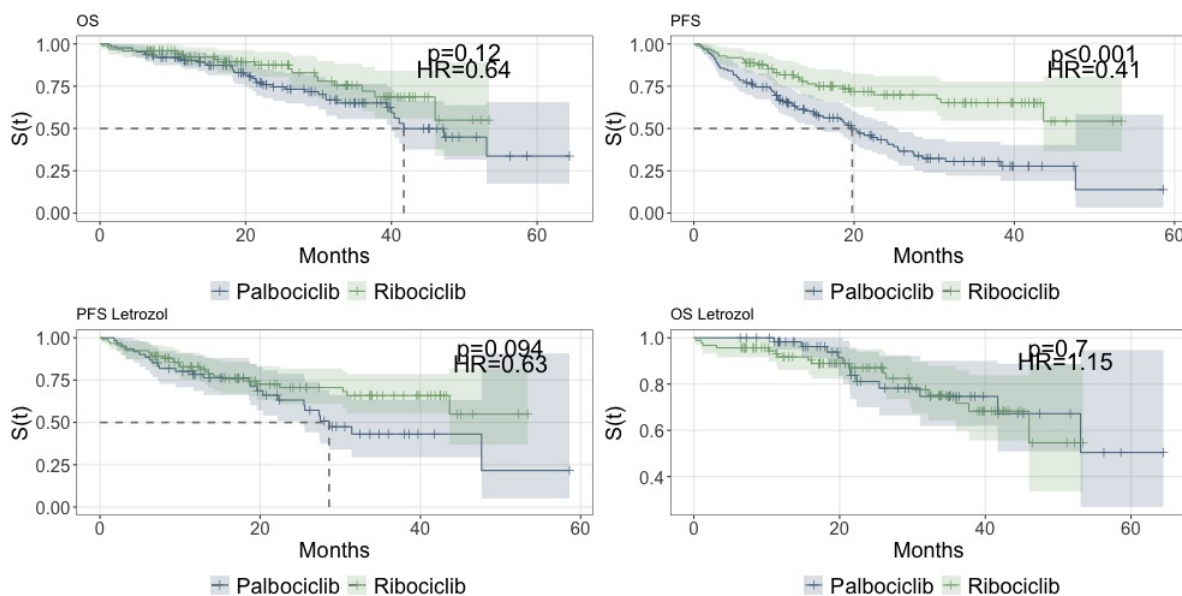
R was used for statistical analysis. Demographic, clinical characteristics and side effects were analyzed using descriptive statistics (count, percentages and median/range). Kaplan–Meier test was used to determine the median PFS and OS in the entire population and subgroups. Log-rank test was used for comparisons of PFS and OS among different subgroups. Cox Regression

was used to assess feature importance and impact. All statistical tests were two-sided, and the significance level was 0.05. The evaluation of the proportional hazards assumptions was done by Schoenfeld residues analysis. We applied propensity score weights to achieve a more robust comparison between the two groups of CDK4/6i. We used the existence of visceral metastases, treatment line, age at treatment start, and stage. We used the `WeightIt` package for R [80]. We applied the weights to the Kaplan-Meier curves and to the Cox Regression. We applied the weights to get the ATE which is  $E[Y_i(1) - Y_i(0)]$ , the average effect of moving an entire population from untreated to treated, or from one drug to the other. Weights were used instead of matching since it is more suited for calculating ATE and the need to preserve the sample size since it is already small from the start. The formula for calculating the weights was through propensity score weighting with GLM. Multiple comparisons were done with the Benjamini-Hochberg (BH) method.

### 3.7.6 Results

The median OS in the entire population treated with CDK4/6 inhibitors was 46 months (95%CI 39.4–55.6). Median PFS was 20.1 months (95%CI 18.3–24.2). Following this, we compared Palbociclib and ribociclib as first-line treatments. We found that regarding OS, there is no significant difference between the two, but ribociclib is significantly better in terms of PFS ( $p\text{-value} \leq 0.001$ ) (figure 3.19). Additionally, we compared the same CDK4/6 inhibitors with letrozol as a combination only. Regarding this scenario, we found out that both were similar in terms of OS and PFS.

Figure 3.19: Survival curves for Palbociclib and Ribociclib (1st line) - Progression Free Survival and Overall Survival



We then compared both with a cox regression, where OS shows no significant difference between palbociclib and ribociclib but a significantly better PFS for ribociclib (figure 3.13)  $HR$  0.60

[95%CI 0.36-0.97] when adjusted to the stage, visceral metastases, age, treatment line, combination and ECOG. This data implies that ribociclib reduces the risk of the disease progression by 40% compared to palbociclib when adjusted for the variables mentioned. The proportional hazards assumption was confirmed with p values all over 0.10.

Table 3.13: Cox Regression with palbociclib and Ribociclib - Progression Free Survival and Overall Survival

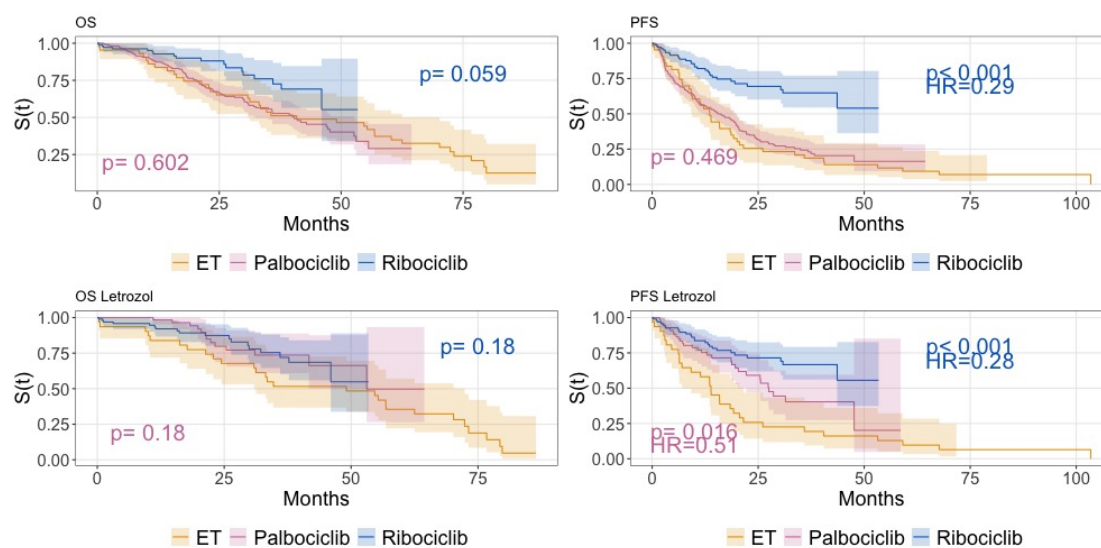
Characteristic	OS			PFS		
	HR <sup>†</sup>	95% CI <sup>†</sup>	p-value	HR <sup>†</sup>	95% CI <sup>†</sup>	p-value
Drug						
Palbociclib	—	—		—	—	
Ribociclib	1.04	0.54, 2.02	>0.9	0.60	0.36, 0.97	0.039
Combination						
Fulvestrant	—	—		—	—	
Letrozol	0.35	0.18, 0.66	0.001	0.45	0.28, 0.71	<0.001
Treatment Line						
1st Line	—	—		—	—	
2nd+ Lines	1.00	0.62, 1.64	>0.9	1.20	0.82, 1.77	0.3
Stage						
I	—	—		—	—	
II	5.67	1.36, 23.6	0.017	1.89	0.99, 3.62	0.055
III	8.17	1.95, 34.2	0.004	3.02	1.57, 5.80	<0.001
IV	8.06	1.91, 34.0	0.004	2.24	1.15, 4.35	0.017
Visceral Metastasis						
No	—	—		—	—	
Yes	1.75	1.19, 2.56	0.004	1.35	1.00, 1.82	0.049
Age at treatment start	1.00	0.98, 1.02	>0.9	0.99	0.97, 1.00	0.053
ECOG at treatment start						
0	—	—		—	—	
1	1.62	1.05, 2.49	0.030	1.22	0.88, 1.69	0.2
2	3.90	2.05, 7.41	<0.001	1.54	0.86, 2.76	0.15
3	12.4	3.55, 43.1	<0.001	0.39	0.05, 2.87	0.4

<sup>†</sup> HR = Hazard Ratio, CI = Confidence Interval

When comparing endocrine therapy with CDK4/6 inhibitors as first-line treatment (figure 3.20), we see that only Ribociclib is significantly better in terms of PFS and OS (p-value  $\leq 0.001$ ). When comparing palbociclib as the first line, we see that there is no significant difference in terms of PFS and OS (p=0.6 and 0.47). We also applied the same analysis as above, comparing only the letrozol combination with letrozol alone. We found that both ribociclib and palbociclib are significantly better in terms of PFS (HR 0.51 for palbociclib and 0.28 for ribociclib).

When comparing palbociclib and ribociclib adjusted for ATE weights, we found a different scenario from previous assessments. There is a significant difference between the two in terms of OS and PFS (figure 3.21). We calculated the weights taking into account stage, age at treatment start, treatment line, and ECOG.

Figure 3.20: Survival curves (OS and PFS) comparing endocrine therapy (ET) to CDK4/6 inhibitors as 1st line. p values shown as pairwise vs. ET.



The Cox regression adjusted for weights shows that ribociclib is not significantly different from palbociclib for OS. The HR for PFS is 0.54 [0.31-0.94;  $p=0.029$ ], implying that ribociclib reduces the risk of the disease progression by ~50% compared to palbociclib when adjusted to the stage, combination drug, treatment line, visceral metastasis, age, and ECOG. Proportional hazard assumptions are confirmed as well.

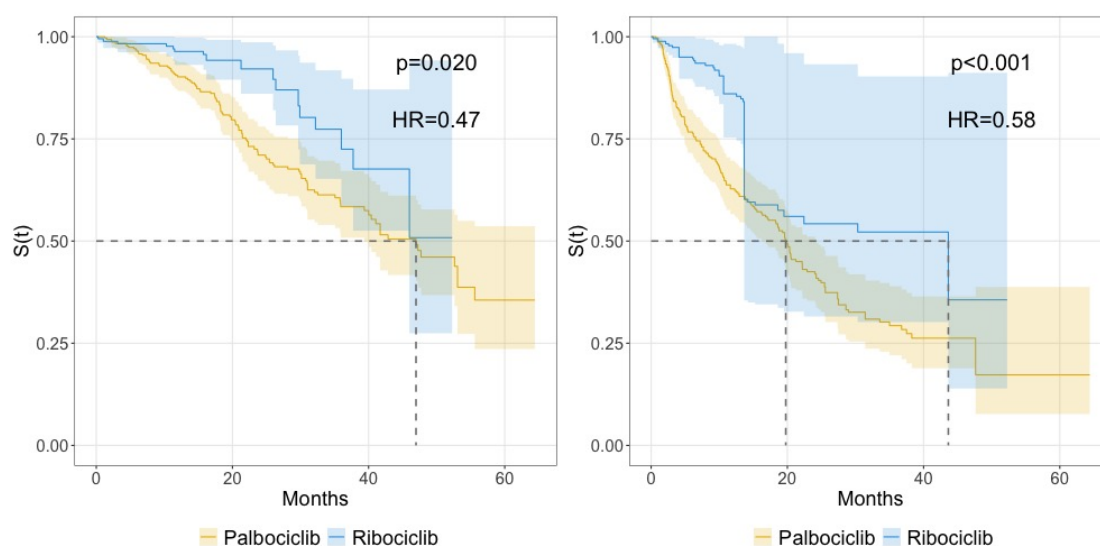
### 3.7.7 Discussion

The aim of this study was to evaluate the real-world use of palbociclib and ribociclib in combination with ET for HR+/HER2– and compare this drug class with traditional endocrine therapy. Few real-world evidence studies of palbociclib and ribociclib used in daily clinical practice have been published identifying clinical benefit, patient profile, and sequencing of treatment, with even less evidence for the Portuguese population.

When comparing with clinical trials, regarding patient profile, in our study, 51% had visceral metastasis and 35% had bone-only metastases compared with 49% and 38% in PALOMA-2, and 60% and 25% in PALOMA-3, respectively [169, 55]. As for ribociclib and bone-only metastases, MONALEESA-7 [191] has 24% and MONALEESA-2 has 40% [88] and our study has 30%.

Of note, the range of median PFS for first-line palbociclib was 15.5–25.5 months, which is shorter than 27.6 months observed in a post hoc analysis of the PALOMA-2 clinical trial with extended follow-up [169], but in line with RWE studies (13.3–20.2 months) [84]. When assessed with only letrozol as a combination, the median PFS increased to 28.6 months [95% CI 25.5–not reached]. As for ribociclib, median survival time was not reached whether in OS and PFS. So we can at least say that the median PFS is longer than 50 months. This is longer than the median progression-free survival of 23.8 months (95% CI 19.2–not reached) reported in the

Figure 3.21: Comparison of palbociclib and ribociclib survival curves adjusted for propensity scores



MONALEESA-7 trial [191] and longer than 25.3 months (95% CI 23.0–30.3) in the MONALEESA-2 trial [88].

Regarding the comparison between ET and CDK4/6i first line, we found out that neither OS and PFS have significant changes when compared ET to Palbociclib 1st line. We can see the values similar to clinical trials when comparing only the letrozol group (both combination and letrozol alone). For this subgroup, we have similar results to clinical trials, with palbociclib being significantly better, with an HR of around 0.5.

Ribociclib is significantly better for the PFS when compared with letrozol and fulvestrant and with letrozol alone, with an HR of around 0.29 for PFS and 0.28 for ribociclib. This would imply that a combination with fulvestrant should be more effective when used with ribociclib and palbociclib. To note, that despite their results, the values in table 3.13 suggest that when we adjust for the variables indicated, ribociclib is significantly better than palbociclib in terms of PFS with an HR of around 0.6.

When comparing with propensity scores weighting, we found out that ribociclib is significantly better than palbociclib for PFS. Our findings suggest that ribociclib could be a better approach for treating HR+, HE- metastatic breast cancer, providing a median OS of over 40 months and median PFS of around 42 months.

### 3.7.8 Conclusion

For conclusions and next steps, we feel we have demonstrated that the ribociclib is a good alternative to palbociclib. We still do not have sufficient evidence to state that palbociclib is actually better than endocrine therapy regarding Overall Survival. However, it is sufficient to state that CDK4/6i have an impact on PFS. Further information about the population could be interesting,



as well as providing information about safety, economic impact, and quality of life. The characterization of the population in terms of biomarkers could be very useful. We aim to address those issues in sequencing papers. Finally, since all of this data was collected from a single institution, we can not generalize the results to the entire population. However, we believe that this study can be used as a starting point for further research in this area. Additionally, this evidence was generated from observational data. Although we adjusted for confounding factors, we cannot exclude the possibility of residual confounding. However, the propensity scores matching allows for a more robust comparison between the two groups, there is still the possibility of unmeasured confounders.

### **3.8 Leveraging data to create Clinical Decision Support Systems**

This section is based on the paper entitled "Machine-learning in Obstetrics: FHIR-based Support System for predicting delivery type". This work was in part a result of the work in section 3.5. While testing for distributed mechanisms, we kind of felt that some evaluation metrics were inspiring to pursue this further. We built a CDSS system that is interoperable and aims to provide support for subpar evaluation of a Cesarean Section (C-Section).

#### **3.8.1 Introduction**

The ability to provide care to both women and newborns during delivery is one of the most important aspects of healthcare and is often used as a metric to assess healthcare as a whole across different countries. C-Sections are one of the most important aspects of delivering babies since it has a considerable impact on the mother's health and well-being. Despite this type of procedure increasing over the last few years, it is still illusive the reasons behind such events. Reports from 2016 suggest that this increment is a global phenomenon, being that from 1990 to 2014, this type of delivery almost increases by 3-fold from 6.7% to 19.1% [27, 42]. Some of these impacts, being more prone to investigation in the last years, including the risk of infection, haemorrhage, organ injury and complications related to the use of anaesthesia or blood transfusion [53, 122]. There is also a higher risk of complications in subsequent pregnancies like uterine rupture, abnormal placental implantation and the need for hysterectomy [104, 185]. As for the infant, C-Sections include the risk of respiratory problems, asthma and obesity in childhood [104]. Facing this, in 2015, World Health Organisation released a statement regarding C-Sections rates. Even when other complications could not be totally assessed, it was concluded that C-Section rates higher than 10% were not associated with a reduction in maternal or newborn mortality [213].

Since there is no evidence that this type of procedure is beneficial for women or babies when there is no clear need for it, the focus on filtering such cases is important [42]. Moreover, particularly in Portugal, C-Sections are used as a way of financing healthcare institutions. This was implemented as a strategy of decreasing C-Sections across the country. A committee was created especially with the purpose of reducing the percentage of C-Sections nationally. One of the actions taken along this creation was the reduction of government funding for hospitals with rates of



C-Sections above 25%. In 2020, the number of C-Sections in Portugal is about 36.3%. Almost at the all-time high of 36.9% in 2009 [149]. So, lowering the proportion of C-Section can provide health and financial benefits to institutions and populations alike. With this in mind, we developed a machine-learning algorithm-based support system to assist clinical teams to detect cases of potentially unnecessary C-Sections for analysis. So in this paper, we propose:

- help to provide a method of bringing to the discussion of clinical staff possible less than optimal care regarding deliveries;
- elaborates on how clinical decision support systems can be developed using interoperability standards;
- understand, based on the gathered data, which are the more impacting features for predicting delivery type outcome;
- open a research path regarding the evaluation of this type of clinical decision support system prior to the delivery;
- Perform a concise economical analysis to assess the potential financial impact of implementing the proposed clinical decision support tool.

### 3.8.2 Rationale and Related Work

Regarding the related work, several teams already tackled the potential of predicting the delivery type before birth. However, we believe that creating such a system may have a huge impact on the clinical team, patient and family regarding expectations. For such a system to be possible to enter clinical practice safely, as we hope this one does, several tests and evaluations should be done before going live. On the other hand, we are aiming for a *post-partum* analysis in order to signalise potential sub-optimal decisions so the clinical teams can evaluate the case afterwards, and hopefully, learn about what could have been done better. Works and studies on this matter have been done before, related to second opinions in the healthcare practice regarding the decision of the C-Section [10] and the implementation of clinical guidelines help as well [36]. We hope to provide support to help teams go in this direction.

Nevertheless, in the literature, there are works related to predicting a successful vaginal birth after a previous C-Section, like the work of Lipschuetz et al., [115] where a gradient boosting method was used to predict such event using prenatal data to do so. Grobman et al., [82] did similar work with a multivariable logistic regression model. There was also the usage of different modalities of data for predicting delivery type. The work of Fergus et al. [70] introduces a method of predicting delivery type using the fetal heart rate signals. Similarly, the work from Saleem et al. [172] proposes a method of predicting delivery type using interactions between fetal heart rate and maternal uterine contraction. Finally, there are also works that focus on predicting delivery mode before childbirth like the work of Ullah et al. [193] where a boosting algorithm was used in order to predict delivery mode with enriched datasets. Also the work of Gimovsky et al. [74] where decision trees were introduced to predict C-Sections by physician group.

However, as far as we know, there was no model tested (even in a controlled setting) in clinical

practice, with no interoperable format of communication like Fast Healthcare Interoperability Resources (FHIR) or employed to *post-partum* setting and finally, none with real-world data related with Portuguese hospital, making that our paper could be a novelty under very different dimensions.

### 3.8.3 Materials

The data was gathered from 9 different Portuguese hospitals regarding obstetric information: data from the mother, several data points about the fetus and delivery mode. The data is from 2019 to 2020. The software for collecting data was the same in every institution and the columns are the same, even though the version of each software differed across hospitals. Across the different hospitals, data rows ranged from 2364 to 18177. The selected variables have the following distributions are shown in table 3.14. The sum of all rows is 73351 rows. The outcome variable had the following distribution as stated in table 3.15

### 3.8.4 Methods

We wrote all of the code in Python 3.9.7 with the usage of the scikit-learn library [146]. All null representations were standardized. Data was prepossessed with the removal of features with high missing rates ( $> 90\%$  overall). All missing value representations were standardized. The imputation process was done using the KNN imputation method (for continuous variables) or a new category (NULLIMP) for categorical variables. For this purpose, the Birth Type was reduced to binary. All assisted birth were merged into vaginal birth and C-Section remained as the other class. Procedures and diagnosis were also used and were encoded as binary features, we took the time to analyse each one of them in order to avoid leakage since there were procedures obviously related to C-Sections and vaginal deliveries. Feature creation was done through the free-text variable relating to the medication prescribed. Features were collected from it and converted into Anatomical Therapeutic Chemical (ATC) Classification Group level 4, which stands for chemical subgroups. We also created some new features from data in the dataset, namely new categories related to the labour and condition of the baby. Also, a few data quality issues were addressed, like impossible values that were transformed into null. In this category, the main issues were Body Mass Index (BMI)/Weight and gestational age. Finally, only a few columns were selected. We used a mixture of surveying the literature and the feature with greater correlation with the outcome. The models tested were Logistic Regression, Decision Tree, Random Forest, 3 different Boosting methods (as implemented by XGBoost, Lightgbm and scikit learn) and a linear model based on Stochastic Gradient Descent. The evaluation was done with repeated stratified cross-validation with 10 splits and 2 repetitions. The API for serving the prediction model was developed with FastAPI.

Finally, a clinical evaluation was carried out with questionnaires sent to several obstetrics specialists in order to assess the validity and possible impact of the model.

Table 3.14: Distribution of feature used for prediction

Variable	M (SD)	Mode [%]
Mother Age	31.0 (5.6)	
Weight pre-pregnancy	65.8 (13.9)	
Weight on admission	78.6 (14.2)	
BMI	25.0 (5.4)	
Previous eutocic delivery	0.4 (0.7)	
Previous vacuum-assisted delivery	0.1 (0.3)	
Previous forceps	0.0 (0.1)	
Previous C-section	0.1 (0.4)	
Fetal presentation on admission		cephalic [26.323 %]
Bishop score	5.5 (3.0)	
Gestational age on admission	38.9 (1.9)	
Premature rupture of the membrane		No [87.991 %]
Chronic hypertension		No [97.676 %]
Gestational hypertension		No [97.749 %]
Preeclampsia		No [98.299 %]
Gestational diabetes		No [89.811 %]
Gestational diabetes treated with diet		No [94.285 %]
Gestational diabetes treated with insulin		No [98.083 %]
Gestational diabetes treated with oral antidiabetic drugs		No [97.797 %]
Maternal Diabetes		No [99.509 %]
Type 1 Diabetes		No [99.816 %]
Type 2 Diabetes		No [99.843 %]
Presentation at birth		Vertex presentation [94.000 %]
Delivery		Spontaneous [53.864 %]
Gestational age on birth	39.0 (1.8)	
Smoking during pregnancy		No [88.442 %]
Alcohol consumption during pregnancy		No [98.65 %]
Consumed drugs during pregnancy		No [99.825 %]
Nr of pregnancies (with current)	1.9 (1.1)	
Pregnancy type		Spontaneous [85.417 %]
Surveillance		yes [97.699 %]
Hospital surveillance		yes [67.807 %]
Pelvis Adequacy		Adequate [17.512 %]
Consistency of the cervix	1.6 (0.6)	
Fetal station	0.8 (0.8)	
Dilation of the cervix	1.3 (0.8)	
Effacement of the cervix	1.2 (1.2)	
Position of the cervix	0.6 (0.7)	
Haematologic disease		No [95.674 %]
Respiratory disease		No [95.605 %]
Cerebral disease		No [98.793 %]
Cardiac disease		No [92.967 %]
Neuroaxis techniques		1 [69.5 %]
Number of children	0.6 (0.8)	

Table 3.15: Distribution of Delivery Methods

Type of delivery	Frequency (%)
C-Section	19 803 (27%)
Vaginal	38 189 (52%)
Instrumental delivery	15 359 (21%)

### 3.8.5 Results

#### 3.8.5.1 The model

The results are mainly the model that predicts the occurrence of a C-Section or natural delivery. The evaluation metrics are present in the table below for the best hyper-parameters found for the training data.

Table 3.16: Performance Metrics in the training set with mean AUROC and 95% Confidence Interval (CI)

Metric	AUC	CI 95%
XGBoost	0.8809	0.8799, 0.882
Decision Tree	0.8337	0.8324, 0.8349
Logistic Regression	0.8716	0.8706, 0.8726
AdaBoost	0.8753	0.874, 0.8766
LightGBM	0.8805	0.8793, 0.8817
Stochastic Gradient Descent	0.8704	0.8694, 0.8713
Random Forest	0.8752	0.8743, 0.8762

XGBoost and LightGBM were the best-performing algorithms. However, we selected LightGBM [103]. since it is faster and requires less memory. The threshold selected for deploying the model was 0.7457247885715557 which rendered the metrics in the test set like it is shown in table 3.17.

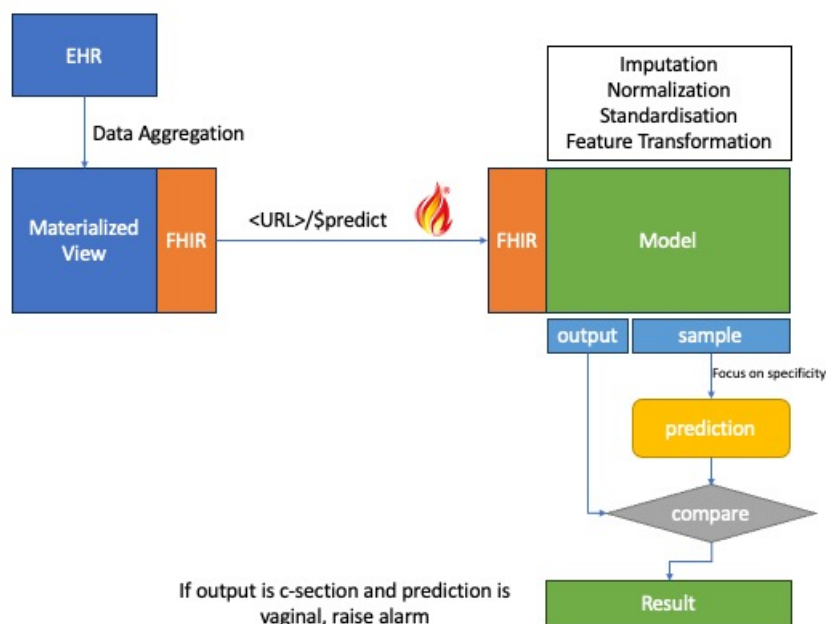
Table 3.17: Performance Metrics in the test set with chosen threshold

Metric	Value
Accuracy	0.8052
Sensitivity	0.8223
Precision	0.9023
F1 Score	0.8605

### 3.8.5.2 Deployment

The purpose of this model is to be served as an API for usage within a healthcare institution and act as a supplementary management decision support tool for obstetrics teams. And for that to happen, a health information system must make the requests to the API. Even though a concrete, vendor-specific information model and input health information system were used, we hope to create a more interoperable clinical decision support system which can be used by every system that acts upon births and obstetrics departments. That is why we built it around the Health Level Seven (HL7) Fast Healthcare Interoperability Resources (FHIR) standard (R5 version) in order to simplify the method of interacting with the API. This decision, opposed as to using a proprietary model for the data, sits upon the usage of FHIR resources: Bundle and Observation for request and returning the result as a message through a custom operation called "\$predict". It is intended to publish the profiles of these objects in order to facilitate access to the API using standardized mechanisms and data models - current build of the profiles <https://joofio.github.io/obs-cdss-fhir/>. The current spec is detailed in the following link. The process is identified in figure 3.22. We deployed this model in production in a single hospital and without a user interface, only collecting the data and prediction for later discussion and analysis. We collected 3231 requests. During this time, the number of alarms that were triggered was 123 (3.8%). From this, we tried to understand the level of certainty for the decision and check the difference from the threshold of these alarms. The distance to the threshold for 73 was lower than 0.1 and was bigger than 0.1 for 50 (1.55%) cases.

Figure 3.22: Deployment and decision mechanism of the model



### 3.8.5.3 Clinical Evaluation

The clinical evaluation was done by sending questionnaires to clinicians with a relationship with obstetrics in order to assess 10 patients, with only access to the variables used by the model and to answer 3 questions for each. The first was to give a score from 1-10 of how likely that patient would give birth through C-Section, then to select the feature/variable that most influenced the decision and which feature they would require to make a better assessment. We sent the questionnaire to 20 people and got 6 answers. This rendered the results in figure 3.23. We also predicted the result with the model as stated in figure 3.23. These patients were new and not seen ever by the model in the training phase. As for the analysis of the most important features and missing features, the missing features were categorized into 3 categories: 1) Existent in the dataset but not included in the model, 2) Non-existent in the dataset and 3) existent in the dataset and included but that particular information was not filled for the patient assessed. This rendered a total of 62% non-existent and 38 % existent but no information was provided at that moment. No feature mentioned existed but had not been included in the model. From the non-existent, 38 % were new clinical assessments, 38% were linked to information from previous births, 15% connected in more in-depth information about provided information (i.e, motive for induction) and 11% were related to the mother's choice (if she wanted a C-Section). As for feature importance, from the 60 answers, we got 55% with labour being the most important factor. 15% answered the number of previous vaginal births, 8% the evolution of weight and another 8% the number of previous C-Sections. The remaining 14% were various features, from BMI, neuroaxis techniques, gestational age and weight of the mother. Of all of these, 90% were included and were in the top 10 features of the model.

### 3.8.5.4 Potential Financial Impact

The financial support provided to public hospitals in Portugal is partially tied to the rate of C-Sections. To assess the potential impact of this mechanism on all Portuguese public hospitals, we conducted a simulation study. We got data for every public hospital for the last 12 months [149] and applied a reduction of 3.8% (the rate of warnings triggered in the new dataset) and recalculated the rate of C-Sections. The increase in support was calculated by the state-mandated rate as shown in table 3.18. With this new rate, we observed that implementing our tool would result in financial benefits for 30% (11 hospitals) of the public hospitals. Specifically, five hospitals would begin receiving support instead of no support at all. Three hospitals would experience a doubling of their financial benefit, while two hospitals would see a 50% increase. Furthermore, one hospital would receive an additional one-third of financial support. If we assumed that only half of the warnings found in the new data were actually true (1.9%) we found that only 6 hospitals would be benefited. 3 from 0 to 0.25, 2 from 0.25 to 0.50 and 1 from 0.50 to 0.75.

Figure 3.23: Validation data. The colour represents the actual birth type. The boxplot represents the median and Inter-Quartile Range of the reviewers and the X represent each patient case. Contains 6 Vaginal births and 4 C-Sections. \* represents wrong predictions of the model.

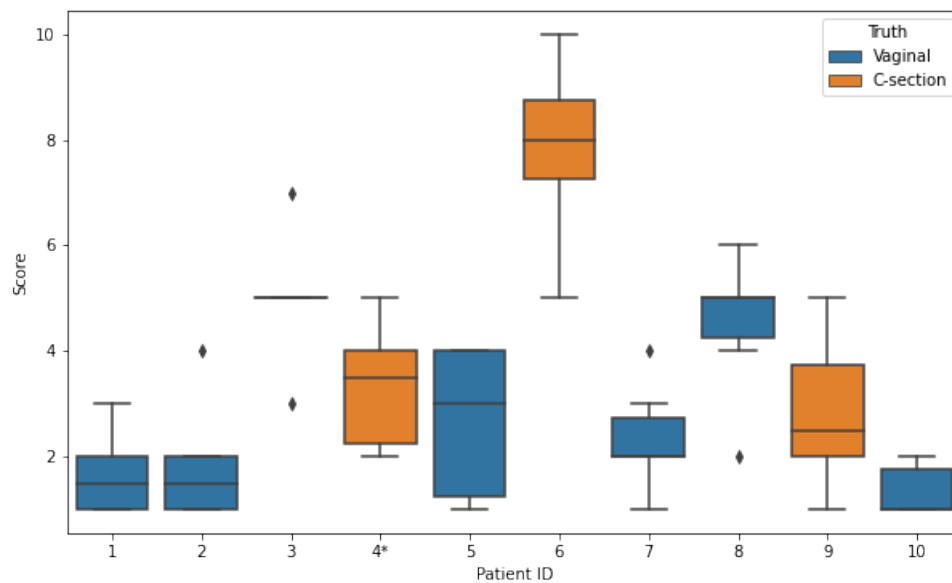


Table 3.18: Ruleset for state-provided financial support indexed to C-Sections. X is the current payment of a C-Section inpatient episode. Adapted from [6]

Rate of C-Sections	Support
<25%	x
[25%, 26.4%]	0.75 x
[26.5%, 27.9%]	0.5 x
[28%, 29.4%]	0.25 x
>29,5%	0

### 3.8.6 Discussion

The first thing to address about this model is the number of biases that we introduced in the model by choice. We joined all vaginal delivery types into a single category (assisted and non-assisted) which introduces a bias since these delivery modes are indeed different. Secondly, the fact that we want to predict if the delivery type was wrongly chosen, mainly for the case of a C-Section that did not need to be so, is also a bias. We used this approach because the initially collected data did not have the representation of such events. So the biases of possibly wrong delivery types were present in the training data. We tried to go around this factor by selecting a threshold that gave the model higher sensitivity than specificity so that only large probabilities would trigger an alarm for human consideration. parallel to this, we are starting to gather labelled cases, with the help of clinicians in order to create a better training dataset. Furthermore, since the data was collected from different hospitals, differences in the data input can also occur. Even though the health information system is the same, the processes that originate the data and are being used for secondary purposes could introduce several biases in the data. This is an issue that was accepted from the start regarding the mechanism of data collection and model training. Regarding the clinical evaluation, it was only possible to get a glimpse of several things due to the number of respondees and the actual model still being implemented at the moment. Despite that, the results are encouraging, since the model seems to behave better than humans with the data provided. However, this is a biased vision, since clinicians in the real world, can have access to more data and information than the model has. It is encouraging, but caution is advised before more testing and evaluation are performed. As for the deployment, further work could be the improvement of the API in order to map all variables to an ontology, making it easier for every system and person to access it and get a suggestion of the delivery type. Finally, as for the next steps, we feel the assessment could still be improved, as for the model, taken that more data is provided and labelling is added as well. One issue interesting as well is the fact that 38% of the answers regarding the most important data element missing from the patient record is data that is being collected but was missing for that patient in specific, raises an important question about data input and data interoperability. If we cannot have access to data when it matters most, it can become meaningless.

### 3.8.7 Conclusion

We believe that we developed a fairly robust system for alarming for possibly wrong C-Sections, that could have a positive impact in real-world practice. However, there are issues to tackle before doing so. There is a need for further evaluating the impact of such a system on clinical practice decisions. There could be a very wide range of reasons that could lead to a possibly sub-par decision regarding delivery type. From the mother's decision to a lack of information at the decisive moment. This system is not meant to create hurdles regarding practice or to point out defective decisions putting certain professionals under the spotlight. All the underlying assumptions and prejudices about having autonomous systems providing support for practice must be taken into account. Nevertheless, the metrics and results so far are definitely encouraging for having a positive



impact on health and economic outcomes.



*We never are definitely right, we can only be sure we are wrong.*

Richard P. Feynman

# 4

## Discussion

Extracting knowledge from healthcare data is not easy. It relies on the availability of data, which is not always the case, and on the ability to extract knowledge from it. In this chapter, we discuss the main challenges we faced during the development of this thesis, and how we overcame them. We also discuss the limitations of our work, and how it can be improved in the future. Finally, we discuss the main contributions of this thesis, and how they can be used to improve the quality of healthcare. The first problem is getting access to data. The data is not always available, and when it is, it is not always in the format we need. Ethics committees and Data Protection Officer (DPO) requirements are put in place in order to guarantee the patient's privacy and security, but a lot of times at the cost of timely access to data. I consider that synthetic data can have a good impact on this work. While we can leave the legal processes be, we may use synthetic data with a heavy focus on security to develop and test our algorithms. This is a very promising area of research, and I believe it will be a game-changer in the future. Parallel to this approach are distributed paradigms. Having a distributed approach to data analysis could be of great help. This would allow for the data to be analyzed in its original location in a more secure way and timely manner. If metrics and models could be built by local teams and shared across regions and/or countries to leverage the power of the many for single institutions could be groundbreaking. However, underlying both these approaches are data dictionaries and data governance tools. Having the correct functional/clinical description of data could be of great impact on the usage of data. Having already the variables defined as categorical, numerical and so on could be of great help. This is a very important aspect of data science, and it is often overlooked. Simple statistics of datasets could be useful as well. For example, the number of missing values, the number of unique values, the number of outliers, and so on. This would help the data scientist to understand the data better

and to know what to expect from it.

This issue also relates to the second big hurdle of knowledge extraction from healthcare data - quality. As discussed in section 3.4, this is a very complex and sometimes elusive concept. In our case, this implied a lot of time spent with data preprocessing. We had to deal with missing values, outliers, and correctness in the context of the records, and data in different formats. We also had to link together different databases from different HISs which brought to light new problems like the new dimensions of correctness of data. There is a common saying that sums this pretty well *When we have one watch, we know the time, but when we have two, we may never know.* So if we had different information regarding the same variable in different systems, how to decide what is true? Another aspect that is often overlooked is the relationship with the clinicians. We need to understand that they are the ones who will use the tools we develop, and they need to be involved in the process. We need to understand their needs and their workflow. We need to understand what they need and how they need it. We need to understand that they are not data scientists, and they do not have the time to learn how to use our tools. We need to make it easy for them to use our tools. Now healthcare is often explained in terms of clinical teams of different backgrounds. A similar concept could be beneficial for harvesting knowledge from data. Thirdly, building software or tools based on this data is still an early subject that possibly requires a legal and technical framework. A legal is connected to the impact of such tools in healthcare. If drugs require such a long time to be approved in order to assess security, how can we approve a tool that can have a similar impact? A technical framework is connected to the fact that we are still in the early stages of a new health data science paradigm. We are still trying to understand how to use data, and how to extract knowledge from it. We are still trying to understand how to evaluate the performance of our tools. We are still trying to understand how to evaluate the impact of our tools in healthcare in a timely manner in a way that is not biased and that is not too expensive. Imposing similar structures to drugs is ill-advised since it could possibly kill the innovation potential and the interest in providing such tools. And this is where a quality infrastructure could be of use. Seriously betting of biomedical informatics could render huge payoffs down the line. Having the human and material resources to build data infrastructures on local (healthcare institutions) and regional, or even country-wise or cross-country policies to use effective use healthcare data is essential. At the time of the writing of this thesis, examples like EHDS are very promising initiatives that could help to overcome the hurdles of data availability and quality. However cross-country initiatives will always be as good as the weakest link, so it is important to have a common framework and a common goal and to have the resources to achieve it. In concrete, having data pipelines, data governance and data interoperability tools, and data quality tools are essential. Having a common data dictionary and a common data format would also be of great help. This would allow for a more efficient use of data, and it would allow for the use of healthcare data to drive innovation. Tightly connected with this is the possibility of having Real World Evidence (RWE) support clinical decisions live. Having data like the one produced in 3.7 in real-time or with high update frequency could be leveraged in order to further support clinicians in making decisions based on data. However, we would require not only the premisses

already discussed, like data quality and cross-collaboration clinics, but a trust-framework would also be necessary. In order to make the automatic dashboard and metrics reliable, transparency is key. Having explainability and transparency in the process of evidence production will be key to building trust and accountability.



*An expert is a person who has made all the mistakes that can be made in a very narrow field.*

Niels Bohr

# 5

## Limitations, future work and conclusions

Regarding limitations, all the projects done in this thesis focus on different aspects of the process of extracting knowledge from healthcare data. Also, they are heavily reliant on specific use cases that are not necessarily generalizable. For example, the 3.7 project is focused on a specific disease, and the 3.8 and 3.5 projects are focused on a specific type of data, and the 3.4 project is focused on a specific type of data and a specific clinical specialty. This means that the results of these projects are not necessarily generalizable to other diseases or other types of data. However, the methods used in these projects are generalizable, and they can be used in other projects. For example, the 3.7 methods can be used to predict in real time, and the 3.8 and 3.5 models can be used to analyze other types of data. The 3.3 method can be used to analyze any type of dataset and incorporated into data pipelines. In future work, I think the groundwork is laid for actually providing assistance to healthcare teams. However, actually deploying real-world CDSSs is seldom an easy task and requires time, money and patience. This is why that part, the actual deployment of the tools, is left for future work. However, we did many tests in the real world and included clinicians in most of our works, so we are confident that the tools are ready to be deployed and create an impact.

For this work to be complete, I had to gather knowledge from different areas. From biology and chemistry, for the healthcare part of it, process design to understand and formalize processes, and math and statistics for machine learning and EDA, interoperability and standards for getting data together, ethics and privacy to gather data with guarantees to the patient's privacy and for creating ethical-aware models. Had to dwell into terminologies and healthcare codification and semantics to interpret data and also get acquainted with some clinical specialties like obstetrics and oncology. Had to collect evidence and make the bridge between RCTs and observational data or Real World Data (RWD) so study design was also needed to bridge the gap. Maybe this is one

of the main issues with this domain, where a different set of skills is required to do everything. The alternative would be a team of different people and honestly, the most successful projects I have seen are the ones that have a team of different people. Finally, the








## A.1 Data Dictionary

Acronym	Description
IA	Mother Age
GS	Blood Group
PI	Weight at the beginning of pregnancy
PAI	Weight on Admission
IMC	BMI
CIG	If Smoker During Pregnancy
APARA	Number of previously born babies
AGESTA	Number of Pregnancies
EA	Number of Previous Eutocic Deliveries with no assistance
VA	Number of Previous Eutocic Deliveries with help of vacuum extraction
FA	Number of Previous Eutocic Deliveries with help of forceps
CA	Number of Previous C-sections
TG	Pregnancy Type (spontaneous, In vitro fertilisation...)
V	If the pregnancy was accompanied by MD
NRCPN	Number of prenatal consultations
VH	If the pregnancy was followed by a MD in a hospital
VP	If the pregnancy was followed by a MD in a private clinic
VCS	If the pregnancy was followed by a MD in a primary care facility
VNH	If the pregnancy was followed by a MD in the same hospital the delivery was made
B	Pelvis Adequacy
AA	Baby's Position on Admission
BS	Bishop Score
BC	Bishop Score Cervical Consistency
BDE	Bishop Score Fetal Station
BDI	Bishop Score Dilatation
BE	Bishop Score Effacement
BP	Bishop Score Cervical Position
IGA	Number of Weeks on Admission
TPEE	If the delivery was spontaneous
TPEI	If the delivery was induced
RPM	If there was a rupture of the amniotic pocket before delivery began
DG	Gestational Diabetes
TP	Delivery Type
ANP	Baby's Position on Delivery
TPNP	Actual Type of Delivery
SGP	Pregnancy Weeks on Delivery
GR	Robson Group



## B.1 C-section assessment questionnaire


**Hospital X**

Mãe	
Idade da grávida no parto	32
Peso da grávida no início da gravidez	45
Peso da grávida quando é internada para ter o bebé.	S/Info.
IMC da grávida no início da gravidez	17.6
Nº partos eutócicos anter.	0
Nº partos distócicos anter. via vaginal com ventosas	0
Nº partos distócicos anter. via vaginal com fórceps	0
Nº de partos cesárianas anteriores	0
Hipertensão Crónica	S/Info.
Hipertensão Gestacional	S/Info.
Hipertensão Pré-eclâmpsia	S/Info.
Diabetes Gestacional	S/Info.
Diabetes Gestacional com Dieta	S/Info.
Diabetes Gestacional com Insulina	S/Info.
Diabetes Gestacional com antidiabéticos orais	S/Info.
Diabetes Materna	S/Info.
Diabetes Tipo 1	S/Info.
Diabetes Tipo 2	S/Info.
Fumou durante a gestação	S/Info.
Ingeriu álcool durante a gestação	S/Info.
Utilizou estupefacientes durante a gestação	S/Info.
Nº de gestações que teve (esta incluída)	2
Tipo de gravidez actual (se foi espontânea, FIV, etc)	Esp.
Se a gravidez foi vigiada (>= 5 consultas)	Sim
Se a gravidez foi vigiada no mesmo hospital do parto	Sim
Nº de filhos nascidos	0
Doença Hematológica	S/Info.
Doença Respiratória	S/Info.
Doença Cerebral	S/Info.
Doença Cardíaca	S/Info.

Parto	
Semanas de gestação na admissão	39.4
Rotura prematura de membranas	S/Info.
Apresentação no parto	Cefálica de vértice
Trabalho de parto	Espontâneo
Semanas de gestação no momento do Parto	39.4
Posição bebé na 1ª verificação no hospital	S/Info.
Avaliação da pelve óssea	S/Info.
BISHOP Consistência	S/Info.
BISHOP Descida	S/Info.
BISHOP Dilatação	S/Info.
BISHOP Extinção	S/Info.
BISHOP Posição	S/Info.
BISHOP Score	S/Info.
Técnicas do neuroeixo	Sim

Histórico Apresentação	
Apresentação semana 39	S/Info.
Apresentação semana 38	S/Info.
Apresentação semana 37	cefálica
Apresentação semana 36	S/Info.
Apresentação semana 35	cefálica
Apresentação semana 34	S/Info.
Apresentação semana 33	S/Info.
Apresentação semana 32	S/Info.
Apresentação semana 31	cefálica
Apresentação semana 30	S/Info.
Apresentação semana 29	S/Info.
Apresentação semana 28	S/Info.
Apresentação semana 27	S/Info.

Evolução Peso	
Percentil peso semana 39	S/Info.
Percentil peso semana 38	S/Info.
Percentil peso semana 37	25th-50th
Percentil peso semana 36	S/Info.
Percentil peso semana 35	10th-25th
Percentil peso semana 34	S/Info.
Percentil peso semana 33	S/Info.
Percentil peso semana 32	S/Info.
Percentil peso semana 31	10th-25th
Percentil peso semana 30	S/Info.
Percentil peso semana 29	S/Info.
Percentil peso semana 28	S/Info.
Percentil peso semana 27	S/Info.

Com base nesta informação, diga de 1 a 10 quão provável seria originar uma cesariana.

Com base nesta informação, qual a característica/variável/elemento apresentado(a) acima que mais impactou a sua decisão?

Qual a característica/variável/elemento não existente que gostaria de ter para avaliar melhor?

## B.2 Data quality questionnaire



Hospital X

### Ficha nº 1

Mãe		Evolução Peso	
Idade da grávida no parto	37.0	Estimativa peso eco 24	S/ Info.
Grupo sanguíneo da grávida	0,RH_POSITIVO	Estimativa peso eco 25	S/ Info.
Peso da grávida no início da gravidez	56.0	Estimativa peso eco 26	S/ Info.
Peso da grávida quando é internada para ter o bebé.	S/ Info.	Estimativa peso eco 27	S/ Info.
IMC da grávida no início da gravidez	21.9	Estimativa peso eco 28	S/ Info.
Nº partos eutócitos anter. via vaginal sem nada	S/ Info.	Estimativa peso eco 29	S/ Info.
Nº partos eutócitos anter. via vaginal com ventosas	S/ Info.	Estimativa peso eco 30	S/ Info.
Nº partos eutócitos anteriores, via vaginal com fórceps	S/ Info.	Estimativa peso eco 31	S/ Info.
Nº de partos cesarianas anteriores	S/ Info.	Estimativa peso eco 32	S/ Info.
Posição bebé na 1ª verificação no hospital	S/ Info.	Estimativa peso eco 33	S/ Info.
BISHOP Score	S/ Info.	Estimativa peso eco 34	2027.0
Semanas de gestação na admissão	34.0	Estimativa peso eco 35	S/ Info.
Hipertensão Crónica	S/ Info.	Estimativa peso eco 36	S/ Info.
Hipertensão Gestacional	S/ Info.	Estimativa peso eco 37	S/ Info.
Hipertensão Pré-eclâmpsia	S/ Info.	Estimativa peso eco 38	S/ Info.
Diabetes Gestacional	S/ Info.	Estimativa peso eco 39	S/ Info.
Diabetes Gestacional com Dieta	S/ Info.	Estimativa peso eco 40	S/ Info.
Diabetes Gestacional com Insulina	S/ Info.	Estimativa peso eco 41	S/ Info.
Diabetes Gestacional com antidiabéticos orais	S/ Info.	Estimativa peso eco 42	S/ Info.
Diabetes Materna	S/ Info.		
Diabetes Tipo 1	S/ Info.		
Diabetes Tipo 2	S/ Info.		
Fumou durante a gestação	S/ Info.		
Ingeriu álcool durante a gestação	S/ Info.		
Utilizou estupefacientes durante a gestação	S/ Info.		
Nº de gestações que teve (esta incluída)	2		
Tipo de gravidez actual (se foi espontânea, FIV, etc)	ESPONTANEA		
Se a gravidez foi vigiada ( >= 5 consultas)	Sim		
Se a gravidez foi vigiada no mesmo hospital do parto	S/ Info.		
Avaliação da pelve óssea	S/ Info.		
Doença Hematológica	S/ Info.		
Doença Respiratória	S/ Info.		
Doença Cerebral	S/ Info.		
Doença Cardíaca	S/ Info.		
Nº de filhos nascidos	S/ Info.		

Parto		Histórico Apresentação	
Tipo de gravidez actual (se foi espontânea, FIV, etc)	ESPONTANEA	Apresentação na semana 42	S/ Info.
Altura uterina. Medição da altura/tamanho da barriga em cm.	S/ Info.	Apresentação na semana 41	S/ Info.
Avaliação da pelve óssea	S/ Info.	Apresentação na semana 40	S/ Info.
Posição bebé na 1ª verificação no hospital	S/ Info.	Apresentação na semana 39	S/ Info.
BISHOP Score	S/ Info.	Apresentação na semana 38	S/ Info.
BISHOP Consistência	S/ Info.	Apresentação na semana 37	S/ Info.
BISHOP Descida	S/ Info.	Apresentação na semana 36	S/ Info.
BISHOP Dilatação	S/ Info.	Apresentação na semana 35	S/ Info.
BISHOP Extinção	S/ Info.	Apresentação na semana 34	cefálica
BISHOP Posição	S/ Info.	Apresentação na semana 33	S/ Info.
Semanas de gestação na admissão	34.0	Apresentação na semana 32	S/ Info.
Indica se o trabalho de parto foi espontâneo	SIM	Apresentação na semana 31	S/ Info.
Indica se o trabalho de parto foi induzido	S/ Info.	Apresentação na semana 30	S/ Info.
Rotura prematura de membranas	S/ Info.	Apresentação na semana 29	S/ Info.
tipo de parto realizado da gravidez actual.	Parto eutócico cefálico	Apresentação na semana 28	S/ Info.
Apresentação no momento do parto	Cefálica de vértice	Apresentação na semana 27	S/ Info.
Trabalho de parto	Espontâneo	Apresentação na semana 26	S/ Info.
Semanas de gestação no momento do Parto	34.1	Apresentação na semana 25	S/ Info.
Classificação de Robson	10	Apresentação na semana 24	S/ Info.



## References

- [1] IEEE standard computer dictionary: A compilation of IEEE standard computer glossaries. pages 1–217.
- [2] *OF THE WISDOM OF THE ANCIENTS*, volume 6 of *Cambridge Library Collection - Philosophy*.
- [3] John graunt on causes of death in the city of london. *Population and Development Review*, 35(2):417–422, 2009.
- [4] Why do 87% of data science projects never make it into production? <https://venturebeat.com/ai/why-do-87-of-data-science-projects-never-make-it-into-production/>, July 2019.
- [5] Sawsan Abdulrahman, Hanine Tout, Hakima Ould-Slimane, A. Mourad, C. Talhi, and M. Guizani. A Survey on Federated Learning: The Journey From Centralized to Distributed On-Site Learning and Beyond. *IEEE Internet of Things Journal*, 2021.
- [6] ACSS. Termos Referência para contratualização de cuidados de saúde no SNS para 2023. [https://www.acss.min-saude.pt/wp-content/uploads/2016/10/Termos-Referencia-Contratualizacao\\_2023.pdf](https://www.acss.min-saude.pt/wp-content/uploads/2016/10/Termos-Referencia-Contratualizacao_2023.pdf), 2023.
- [7] Amina Adadi and Mohammed Berrada. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access*, 6:52138–52160, 2018.
- [8] Julia Adler-Milstein and Ashish K. Jha. HITECH Act Drove Large Gains In Hospital Electronic Health Record Adoption. *Health Affairs*, 36(8):1416–1422, August 2017. <http://www.healthaffairs.org/doi/10.1377/hlthaff.2016.1651>.
- [9] Jan Philipp Albrecht. How the GDPR Will Change the World. *European Data Protection Law Review*, 2(3):287–289, 2016. Publisher: Lexxion Publisher; <https://web.archive.org/web/20211014090922/https://edpl.lexxion.eu/article/EDPL/2016/3/4>.
- [10] Fernando Althabe, José M. Belizán, José Villar, Sophie Alexander, Eduardo Bergel, Silvina Ramos, Mariana Romero, Allan Donner, Gunilla Lindmark, Ana Langer, Ubaldo Farnot, José G. Cecatti, Guillermo Carroli, and Edgar Kestler. Mandatory second opinion to reduce rates of unnecessary caesarean sections in Latin America: A cluster randomised controlled trial. *The Lancet*, 363(9425):1934–1940, 2004-06-12.
- [11] Roberto Álvarez Sánchez, Andoni Beristain Iraola, Gorka Epelde Unanue, and Paul Carlin. TAQIH, a tool for tabular data quality assessment and improvement in the context of health data. *Computer Methods and Programs in Biomedicine*, 181:104824, November 2019.
- [12] Jessica S Ancker, Lisa M Kern, Alison Edwards, Sarah Nosal, Daniel M Stein, Diane Hauser, and Rainu Kaushal. How is the electronic health record being used? Use of EHR data to assess physician-level variability in technology use. 21(6):1001–1008.

- [13] Ankur Ankan and Abinash Panda. pgmpy: Probabilistic graphical models using python. In *Proceedings of the 14th Python in Science Conference (SCIPY 2015)*. Citeseer, 2015.
- [14] Peter C. Austin. An Introduction to Propensity Score Methods for Reducing the Effects of Confounding in Observational Studies. 46(3):399–424.
- [15] Peter C. Austin. An Introduction to Propensity Score Methods for Reducing the Effects of Confounding in Observational Studies. 46(3):399–424.
- [16] Peter C. Austin. The performance of different propensity-score methods for estimating differences in proportions (risk differences or absolute risk reductions) in observational studies. 29(20):2137–2148.
- [17] Peter C Austin. The use of propensity score methods with survival or time-to-event outcomes: Reporting measures of effect similar to those used in randomized experiments. *Statistics in Medicine*, 33(7):1242–1258, 2014.
- [18] Mrinal Kanti Baowaly. medWGAN Repository. <https://github.com/baowaly/SynthEHR>.
- [19] Mrinal Kanti Baowaly, Chia Ching Lin, Chao Lin Liu, and Kuan Ta Chen. Synthesizing electronic health records using improved generative adversarial networks. *Journal of the American Medical Informatics Association*, 26(3):228–241, 2019.
- [20] Mrinal Kanti Baowaly, Chia Ching Lin, Chao Lin Liu, and Kuan Ta Chen. Synthesizing electronic health records using improved generative adversarial networks. *Journal of the American Medical Informatics Association*, 26(3):228–241, 2019.
- [21] Mrinal Kanti Baowaly, Chao-Lin Liu, and Kuan-Ta Chen. Realistic Data Synthesis Using Enhanced Generative Adversarial Networks. In *2019 IEEE SECOND INTERNATIONAL CONFERENCE ON ARTIFICIAL INTELLIGENCE AND KNOWLEDGE ENGINEERING (AIKE)*, pages 289–292. IEEE; IEEE Comp Soc, 2019.
- [22] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning: Limitations and Opportunities*. fairmlbook.org, 2019. <http://www.fairmlbook.org>.
- [23] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador Garcia, Sergio Gil-Lopez, Daniel Molina, Richard Benjamins, Raja Chatila, and Francisco Herrera. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58:82–115, 2020-06.
- [24] Anna Okula Basile and Marylyn DeRiggi Ritchie. Informatics and machine learning to define the phenotype. *Expert Review of Molecular Diagnostics*, 18(3):219–226, 2018-03-04.
- [25] Brett Beaulieu-Jones, Casey Greene, and Steven Wu. SPRINT-GAN Repository. [https://github.com/greenelab/SPRINT\\_gan](https://github.com/greenelab/SPRINT_gan).
- [26] Brett K. Beaulieu-Jones, Zhiwei Steven Wu, Chris Williams, Ran Lee, Sanjeev P. Bhavnani, James Brian Byrd, and Casey S. Greene. Privacy-Preserving Generative Deep Neural Networks Support Clinical Data Sharing. *Circulation: Cardiovascular Quality and Outcomes*, 12(7):139–148, July 2019.



- [27] Ana Pilar Betrán, Jianfeng Ye, Anne-Beth Moller, Jun Zhang, A. Metin Gülmezoglu, and Maria Regina Torloni. The Increasing Trend in Caesarean Section Rates: Global, Regional and National Estimates: 1990-2014. *PLoS ONE*, 11(2):e0148343, 2016-02-05.
- [28] Jiang Bian, Tianchen Lyu, Alexander Loiacono, Tonatiuh Mendoza Viramontes, Gloria Lipori, Yi Guo, Yonghui Wu, Mattia Prosperi, Thomas J. George, Christopher A. Harle, Elizabeth A. Shenkman, and William Hogan. Assessing the practice of data quality evaluation in a national clinical data research network through a systematic scoping review in the era of real-world data. *Journal of the American Medical Informatics Association: JAMIA*, 27(12):1999–2010, December 2020.
- [29] Vadim Borisov, Tobias Leemann, Kathrin Seßler, Johannes Haug, Martin Pawelczyk, and Gjergji Kasneci. Deep Neural Networks and Tabular Data: A Survey, June 2022.
- [30] Vadim Borisov, Tobias Leemann, Kathrin Seßler, Johannes Haug, Martin Pawelczyk, and Gjergji Kasneci. Deep Neural Networks and Tabular Data: A Survey. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–21, 2022.
- [31] Bauke Brenninkmeijer. WGAN-DP Repository. <https://github.com/Baukebrennkmeijer/On-the-Generation-and-Evaluation-of-Synthetic-Tabular-Data-using-GANs>.
- [32] Bauke Brenninkmeijer. *On the Generation and Evaluation of Tabular Data using GANs*. PhD thesis, Radboud University, 2019.
- [33] Stephen Burgess and Simon G. Thompson. Improving bias and coverage in instrumental variable analysis with weak instruments for continuous and binary outcomes. 31(15):1582–1600.
- [34] Ramiro Camino. mc-medGAN Repository. <https://github.com/rcamino/multi-categorical-gans>.
- [35] Ramiro Camino, Christian Hammerschmidt, and Radu State. Generating Multi-Categorical Samples with Generative Adversarial Networks. *ArXiv*, 2018. arXiv: 1807.01202.
- [36] Nils Chaillet, Alexandre Dumont, Michal Abrahamowicz, Jean-Charles Pasquier, Francois Audibert, Patricia Monnier, Haim A. Abenhaim, Eric Dubé, Marylène Dugas, Rebecca Burne, and William D. Fraser. A Cluster-Randomized Trial to Reduce Cesarean Delivery Rates in Quebec. *New England Journal of Medicine*, 372(18):1710–1721, 2015-04-30.
- [37] Kunal Chandiramani, Dhruv Garg, and N Maheswari. Performance analysis of distributed and federated learning models on private data. *Procedia Computer Science*, 165:349–355, 2019. 11.
- [38] Pete Chapman, Julian Clinton, Randy Kerber, Thomas Khabaza, Thomas Reinartz, Colin Shearer, and Rudiger Wirth. Crisp-dm 1.0 step-by-step data mining guide. Technical report, The CRISP-DM consortium, August 2000.
- [39] Danton S. Char, Nigam H. Shah, and David Magnus. Implementing Machine Learning in Health Care — Addressing Ethical Challenges. *The New England journal of medicine*, 378(11):981–983, March 2018.

- [40] Nitesh V. Chawla, Kevin W. Bowyer, Lawrence O. Hall, and W. Philip Kegelmeyer. SMOTE: synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16(1):321–357, 2002.
- [41] Changyao Chen. Rank-biased Overlap (RBO). [https://github.com/changyaocen/rbo](https://github.com/changyaochen/rbo), 2023-03-14T02:55:37Z.
- [42] Innle Chen, Newton Opiyo, Emma Tavender, Sameh Mortazhejri, Tamara Rader, Jennifer Petkovic, Sharlini Yogasingam, Monica Taljaard, Sugandha Agarwal, Malinee Laopaiboon, Jason Wasiak, Suthit Khunpradit, Pisake Lumbiganon, Russell L Gruen, and Ana Pilar Betran. Non-clinical interventions for reducing unnecessary caesarean section. *The Cochrane Database of Systematic Reviews*, 2018(9):CD005528, 2018-09-28.
- [43] Travers Ching, Daniel S. Himmelstein, Brett K. Beaulieu-Jones, Alexandr A. Kalinin, Brian T. Do, Gregory P. Way, Enrico Ferrero, Paul-Michael Agapow, Michael Zietz, Michael M. Hoffman, Wei Xie, Gail L. Rosen, Benjamin J. Lengerich, Johnny Israeli, Jack Lanchantin, Stephen Woloszynek, Anne E. Carpenter, Avanti Shrikumar, Jinbo Xu, Evan M. Cofer, Christopher A. Lavender, Srinivas C. Turaga, Amr M. Alexandari, Zhiyong Lu, David J. Harris, Dave DeCaprio, Yanjun Qi, Anshul Kundaje, Yifan Peng, Laura K. Wiley, Marwin H. S. Segler, Simina M. Boca, S. Joshua Swamidass, Austin Huang, Anthony Gitter, and Casey S. Greene. Opportunities and obstacles for deep learning in biology and medicine. *Journal of The Royal Society Interface*, 15(141):20170387, 2018-04-04.
- [44] Edward Choi. medGAN Repository. <https://github.com/mp2893/medgan>.
- [45] Edward Choi, Siddharth Biswal, Bradley Malin, Jon Duke, Walter F. Stewart, and Jimeng Sun. Generating Multi-label Discrete Patient Records using Generative Adversarial Networks, 2017. arXiv: 1703.06490.
- [46] Edward Choi, Siddharth Biswal, Bradley Malin, Jon Duke, Walter F. Stewart, and Jimeng Sun. Generating Multi-label Discrete Patient Records using Generative Adversarial Networks, 2017.
- [47] Vikram S. Chundawat, Ayush K. Tarun, Murari Mandal, Mukund Lahoti, and Pratik Narang. TabSynDex: A Universal Metric for Robust Evaluation of Synthetic Tabular Data, 2022-07-12.
- [48] Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46, 1960.
- [49] Herkulaas MvE Combrink, Vukosi Marivate, and Benjamin Rosman. Comparing Synthetic Tabular Data Generation Between a Probabilistic Model and a Deep Learning Model for Education Use Cases, 2022-10-16.
- [50] Comissão Nacional Proteção de dados. Princípios aplicáveis aos tratamentos de dados efetuados no âmbito da investigação clínica, 2015.
- [51] European Commission. A definition of AI: Main Capabilities and Disciplines. Technical report, European Commission, 2019. <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>.
- [52] European Commission. Ethics Guidelines For Trustworthy AI. Technical report, European Commission, 2019. <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>.

- [53] J. R. Cook, S. Jarvis, M. Knight, and M. K. Dhanjal. Multiple repeat caesarean section in the UK: Incidence and consequences to mother and child. A national, prospective, cohort study. *BJOG: an international journal of obstetrics and gynaecology*, 120(1):85–91, 2013-01.
- [54] Kristin M. Corey, Joshua Helmkamp, Morgan Simons, Lesley Curtis, Keith Marsolo, Suresh Balu, Michael Gao, Marshall Nichols, Joshua Watson, Leila Mureebe, Allan D. Kirk, and Mark Sendak. Assessing Quality of Surgical Real-World Data from an Automated Electronic Health Record Pipeline. *Journal of the American College of Surgeons*, 230(3):295–305.e12, March 2020.
- [55] Massimo Cristofanilli, Nicholas C. Turner, Igor Bondarenko, Jungsil Ro, Seock-Ah Im, Norikazu Masuda, Marco Colleoni, Angela DeMichele, Sherene Loi, Sunil Verma, Hiroji Iwata, Nadia Harbeck, Ke Zhang, Kathy Puyana Theall, Yuqiu Jiang, Cynthia Huang Bartlett, Maria Koehler, and Dennis Slamon. Fulvestrant plus palbociclib versus fulvestrant plus placebo for treatment of hormone-receptor-positive, HER2-negative metastatic breast cancer that progressed on previous endocrine therapy (PALOMA-3): Final analysis of the multicentre, double-blind, phase 3 randomised controlled trial. 17(4):425–439.
- [56] Ricardo Cruz-Correia, Pedro Rodrigues, Alberto Freitas, Filipa Almeida, Rong Chen, and Altamiro Costa-Pereira. Data Quality and Integration Issues in Electronic Health Records. In Vagelis Hristidis, editor, *Information Discovery on Electronic Health Records*, volume 12. Chapman and Hall/CRC.
- [57] Sabyasachi Dash, Sushil Kumar Shakyawar, Mohit Sharma, and Sandeep Kaushik. Big data in healthcare: Management, analysis and future prospects. 6(1):54.
- [58] Clinical Practice Research Datalink. CPRD cardiovascular disease synthetic dataset. Medium: Text/CSV Version Number: 2020.06.001 Type: dataset.
- [59] Neil M. Davies, George Davey Smith, Frank Windmeijer, and Richard M. Martin. Issues in the reporting and conduct of instrumental variable studies: A systematic review. 24(3):363–369.
- [60] Nelis J. de Vos. kmodes categorical clustering library. <https://github.com/nicodv/kmodes>, 2015–2021.
- [61] Timo M. Deist, A. Jochems, Johan van Soest, Georgi Nalbantov, Cary Oberije, Seán Walsh, Michael Eble, Paul Bulens, Philippe Coucke, Wim Dries, Andre Dekker, and Philippe Lambin. Infrastructure and distributed learning methodology for privacy-preserving multi-centric rapid learning health care: euroCAT. *Clinical and Translational Radiation Oncology*, 4:24–31, May 2017.
- [62] Barbara Di Camillo, Giuseppe Nicosia, Francesca Buffa, and Benny Lo. Guest Editorial Data Science in Smart Healthcare: Challenges and Opportunities. 24(11):3041–3043.
- [63] Flavio Di Martino and Franca Delmastro. Explainable AI for clinical and remote health applications: A survey on tabular and time series data. *Artificial Intelligence Review*, pages 1–55, 2022-10-26.
- [64] DPautoGAN. DPAutoGAN Repository. <https://github.com/DPautoGAN/DPautoGAN>.

- [65] Anurag Dutt. SMOOTH-GAN Repository. [https://github.com/anuragdutt/synthehr\\_medgan](https://github.com/anuragdutt/synthehr_medgan).
- [66] Khaled El Emam, Elizabeth Jonker, Luk Arbuckle, and Bradley Malin. A systematic review of re-identification attacks on health data. *PLoS ONE*, 6(12), 2011.
- [67] Khaled El Emam, Lucy Mosquera, Xi Fang, and Alaa El-Hussuna. Utility Metrics for Evaluating Synthetic Health Data Generation Methods: Validation Study. *JMIR Medical Informatics*, 10(4):e35734, April 2022.
- [68] European Commission. Proposal for a regulation of the european parliament and of the council on the european health data space, 2022. [https://eur-lex.europa.eu/resource.html?uri=cellar:dbfd8974-cb79-11ec-b6f4-01aa75ed71a1.0001.02/DOC\\_1&format=PDF](https://eur-lex.europa.eu/resource.html?uri=cellar:dbfd8974-cb79-11ec-b6f4-01aa75ed71a1.0001.02/DOC_1&format=PDF).
- [69] Usama Fayyad, Gregory Piatetsky-Shapiro, and Padhraic Smyth. From data mining to knowledge discovery in databases. *AI Magazine*, 17(3):37, March 1996. <https://ojs.aaai.org/aimagazine/index.php/aimagazine/article/view/1230>.
- [70] Paul Fergus, Abir Hussain, Dhiya Al-Jumeily, De-Shuang Huang, and Nizar Bouguila. Classification of caesarean section and normal vaginal deliveries using foetal heart rate signals and advanced machine learning algorithms. *BioMedical Engineering OnLine*, 16(1):89, 2017.
- [71] Richard S Finn, John P Crown, Istvan Lang, Katalin Boer, Igor M Bondarenko, Sergey O Kulyk, Johannes Ettl, Ravindranath Patel, Tamas Pinter, Marcus Schmidt, Yaroslav Shparyk, Anu R Thummala, Nataliya L Voytko, Camilla Fowst, Xin Huang, Sindy T Kim, Sophia Randolph, and Dennis J Slamon. The cyclin-dependent kinase 4/6 inhibitor palbociclib in combination with letrozole versus letrozole alone as first-line treatment of oestrogen receptor-positive, HER2-negative, advanced breast cancer (PALOMA-1/TRIO-18): A randomised phase 2 study. *The Lancet Oncology*, 16(1):25–35, January 2015. <https://linkinghub.elsevier.com/retrieve/pii/S1470204514711593>.
- [72] Richard S. Forsyth. Liver Disorders. UCI Machine Learning Repository, 1990. DOI: [10.24432/C54G67](https://doi.org/10.24432/C54G67).
- [73] Mark J. Giganti, Bryan E. Shepherd, Yanink Caro-Vega, Paula M. Luz, Peter F. Rebeiro, Marcelle Maia, Gaetane Julmiste, Claudia Cortes, Catherine C. McGowan, and Stephany N. Duda. The impact of data quality and source data verification on epidemiologic inference: A practical application using HIV observational data. *BMC public health*, 19(1):1748, December 2019.
- [74] Alexis C. Gimovsky, Daisy Zhuo, Jordan T. Levine, Jack Dunn, Maxime Amarm, and Alan M. Peaceman. Benchmarking cesarean delivery rates using machine learning-derived optimal classification trees. *Health Services Research*, n/a, 2021. [\\_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1111/1475-6773.13921](https://onlinelibrary.wiley.com/doi/pdf/10.1111/1475-6773.13921).
- [75] Matthew P. Goetz, Masakazu Toi, Mario Campone, Joohyuk Sohn, Shani Paluch-Shimon, Jens Huober, In Hae Park, Olivier Trédan, Shin-Cheh Chen, Luis Manso, Orit C. Freedman, Georgina Garnica Jaliffe, Tammy Forrester, Martin Frenzel, Susana Barriga, Ian C. Smith, Nawel Bourayou, and Angelo Di Leo. MONARCH 3: Abemaciclib As Initial Therapy for Advanced Breast Cancer. *Journal of Clinical Oncology: Official Journal of the American Society of Clinical Oncology*, 35(32):3638–3646, November 2017.

- [76] Andre Goncalves, Priyadip Ray, Braden Soper, Jennifer Stevens, Linda Coyle, and Ana Paula Sales. Generation and evaluation of synthetic patient data. *BMC Medical Research Methodology*, 20(1):1–40, 2020. Publisher: BMC Medical Research Methodology.
- [77] Andre Goncalves, Priyadip Ray, Braden Soper, Jennifer Stevens, Linda Coyle, and Ana Paula Sales. Generation and evaluation of synthetic patient data. *BMC Medical Research Methodology*, 20(1):1–40, 2020.
- [78] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Networks. *Communications of the ACM*, 63(11):139–144, June 2014. arXiv: 1406.2661.
- [79] Trisha Greenhalgh. *How to Read a Paper: The Basics of Evidence-based Medicine and Healthcare*. How To. Wiley, 6 edition.
- [80] Noah Greifer. *WeightIt: Weighting for Covariate Balance in Observational Studies*, 2023. <https://ngreifer.github.io/WeightIt/>, <https://github.com/ngreifer/WeightIt>.
- [81] Léo Grinsztajn, Edouard Oyallon, and Gaël Varoquaux. Why do tree-based models still outperform deep learning on tabular data?, July 2022.
- [82] William A. Grobman, Yinglei Lai, Mark B. Landon, Catherine Y. Spong, Kenneth J. Leveno, Dwight J. Rouse, Michael W. Varner, Atef H. Moawad, Steve N. Caritis, Margaret Harper, Ronald J. Wapner, Yoram Sorokin, Menachem Miodovnik, Marshall Carpenter, Mary J. O’Sullivan, Baha M. Sibai, Oded Langer, John M. Thorp, Susan M. Ramin, Brian M. Mercer, and National Institute of Child Health and Human Development (NICHD) Maternal-Fetal Medicine Units Network (MFMU). Development of a nomogram for prediction of vaginal birth after cesarean delivery. *Obstetrics and Gynecology*, 109(4):806–812, 2007.
- [83] R. W. Hamming. Error detecting and error correcting codes. *The Bell System Technical Journal*, 29(2):147–160, 1950.
- [84] Nadia Harbeck, Meaghan Bartlett, Dean Spurdén, Becky Hooper, Lin Zhan, Emily Rosta, Chris Cameron, Debanjali Mitra, and Anna Zhou. CDK4/6 inhibitors in HR+/HER2- advanced/metastatic breast cancer: A systematic literature review of real-world evidence studies. 17(16):2107–2122.
- [85] healthdatainsight.org.uk. The simulacrum. <https://healthdatainsight.org.uk/project/the-simulacrum/>.
- [86] JaWanna Henry, Yuriy Pylypchuk, Talisha Searcy, and Vaishali Patel. Adoption of Electronic Health Record Systems among U.S. Non-Federal Acute Care Hospitals: 2008-2015. Technical report, ONC, 2016.
- [87] M. A. Hernán. A definition of causal effect for epidemiological research. 58(4):265–271.
- [88] G. Hortobagyi, S. Stemmer, H. Burris, Y. Yap, G. Sonke, S. Paluch-Shimon, M. Campone, K. Petráková, K. Blackwell, E. Winer, W. Janni, S. Verma, P. Conte, C. Arteaga, D. Cameron, S. Mondal, F. Su, M. Miller, M. Elmeliy, C. Germa, and J. O’Shaughnessy. Updated results from MONALEESA-2, a phase III trial of first-line ribociclib plus letrozole versus placebo plus letrozole in hormone receptor-positive, HER2-negative advanced breast cancer. *Annals of oncology : official journal of the European Society for Medical Oncology*, 2018.

- [89] George Hripcsak, Jon D. Duke, Nigam H. Shah, Christian G. Reich, Vojtech Huser, Martijn J. Schuemie, Marc A. Suchard, Rae Woong Park, Ian Chi Kei Wong, Peter R. Rijnbeek, Johan van der Lei, Nicole Pratt, G. Niklas Norén, Yu-Chuan Li, Paul E. Stang, David Madigan, and Patrick B. Ryan. Observational Health Data Sciences and Informatics (OHDSI): Opportunities for Observational Researchers. *Studies in Health Technology and Informatics*, 216:574–578, 2015.
- [90] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of Classification*, 2(1):193–218, 1985-12-01.
- [91] Tim Hulsen. Sharing Is Caring—Data Sharing Initiatives in Healthcare. *International Journal of Environmental Research and Public Health*, 17(9):3046, 2020-01.
- [92] Vojtech Huser, Frank J. DeFalco, Martijn Schuemie, Patrick B. Ryan, Ning Shang, Mark Velez, Rae Woong Park, Richard D. Boyce, Jon Duke, Ritu Khare, Levon Utidjian, and Charles Bailey. Multisite Evaluation of a Data Quality Tool for Patient-Level Clinical Data Sets. *EGEMS (Washington, DC)*, 4(1):1239, 2016.
- [93] integraal kankercentrum Nederland. Synthetische dataset NKR beschikbaar voor onderzoekers. <https://iknl.nl/nieuws/2021/synthetische-data-nkr-beschikbaar-voor-onderzoeker>.
- [94] Stefanie James, Chris Harbron, Janice Branson, and Mimmi Sundler. Synthetic data use: exploring use cases to optimise data utility. *Discover Artificial Intelligence*, 1(1):15, 2021.
- [95] Andras Janosi, William Steinbrunn, Matthias Pfisterer, and Robert Detrano. Heart Disease. UCI Machine Learning Repository, 1988. DOI: [10.24432/C52P4X](https://doi.org/10.24432/C52P4X).
- [96] Divya Jatain, Vikram Singh, and Naveen Dahiya. A contemplative perspective on federated machine learning: Taxonomy, threats & vulnerability assessment and challenges. *Journal of King Saud University - Computer and Information Sciences*, 2021-06-05.
- [97] Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3:160035, 2016.
- [98] James Jordon, Jinsung Yoon, and Mihaela Van Der Schaar. PATE-GaN: Generating synthetic data with differential privacy guarantees. *7th International Conference on Learning Representations, ICLR 2019*, pages 1–21, 2019.
- [99] Erik Joukes, Nicolette F. de Keizer, Martine C. de Bruijne, Ameen Abu-Hanna, and Ronald Cornet. Impact of Electronic versus Paper-Based Recording before EHR Implementation on Health Care Professionals’ Perceptions of EHR Use, Data Quality, and Data Reuse. *Applied Clinical Informatics*, 10(2):199–209, March 2019.
- [100] Michael G. Kahn, Tiffany J. Callahan, Juliana Barnard, Alan E. Bauck, Jeff Brown, Bruce N. Davidson, Hossein Estiri, Carsten Goerg, Erin Holve, Steven G. Johnson, Siaw-Teng Liaw, Marianne Hamilton-Lopez, Daniella Meeker, Toan C. Ong, Patrick Ryan, Ning Shang, Nicole G. Weiskopf, Chunhua Weng, Meredith N. Zozus, and Lisa Schilling. A Harmonized Data Quality Assessment Terminology and Framework for the Secondary Use of Electronic Health Record Data. *eGEMS*, 4(1):1244, September 2016.



- [101] U. Kamath and J. Liu. *Explainable Artificial Intelligence: An Introduction to Interpretable Machine Learning*. Springer International Publishing, 2021.
- [102] Lorenz A. Kapsner, Jonathan M. Mang, Sebastian Mate, Susanne A. Seuchter, Abishaa Vengadeswaran, Franziska Bathelt, Noemi Deppenwiese, Dennis Kadioglu, Detlef Kraska, and Hans-Ulrich Prokosch. Linking a Consortium-Wide Data Quality Assessment Tool with the MIRACUM Metadata Repository. *Applied Clinical Informatics*, 12(4):826–835, August 2021.
- [103] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems*, 30:3146–3154, 2017.
- [104] Oonagh E. Keag, Jane E. Norman, and Sarah J. Stock. Long-term risks and benefits associated with cesarean delivery for mother, baby, and subsequent pregnancies: Systematic review and meta-analysis. *PLoS medicine*, 15(1):e1002494, 2018.
- [105] M. G. Kendall. The treatment of ties in ranking problems. *Biometrika*, 33:239–251, 1945.
- [106] E. Kilsdonk, L. W. Peute, and M. W. M. Jaspers. Factors influencing implementation success of guideline-based clinical decision support systems: A systematic review and gaps analysis. *International Journal of Medical Informatics*, 98:56–64, February 2017.
- [107] Oren Kramer, Adir Even, Idit Matot, Yohai Steinberg, and Yuval Bitan. The impact of data quality defects on clinical decision-making in the intensive care unit. *Computer Methods and Programs in Biomedicine*, 209:106359, September 2021.
- [108] Clemens Scott Kruse, Anna Stein, Heather Thomas, and Harmander Kaur. The use of Electronic Health Records to Support Population Health: A Systematic Review of the Literature. *Journal of Medical Systems*, 42(11):214, November 2018. <http://link.springer.com/10.1007/s10916-018-1075-6>.
- [109] Matt J. Kusner and José Miguel Hernández-Lobato. GANS for Sequences of Discrete Elements with the Gumbel-softmax Distribution, 2016. arXiv: 1611.04051.
- [110] Illidan Lab. dp-GAN Repository. <https://github.com/illidanlab/dpgan>.
- [111] Geun Hyeong Lee and Soo-Yong Shin. Federated learning on clinical benchmark data: Performance assessment. *Journal of Medical Internet Research*, 22(10), 2020. 9.
- [112] Geun Hyeong Lee and Soo-Yong Shin. Federated Learning on Clinical Benchmark Data: Performance Assessment. *Journal of Medical Internet Research*, 22(10), 2020-10-26.
- [113] Siaw-Teng Liaw, Jason Guan Nan Guo, Sameera Ansari, Jitendra Jonnagaddala, Myron Anthony Godinho, Alder Jose Borelli, Simon de Lusignan, Daniel Capurro, Harshana Liyanage, Navreet Bhattal, Vicki Bennett, Jaclyn Chan, and Michael G. Kahn. Quality assessment of real-world data repositories across the data life cycle: A literature review. *Journal of the American Medical Informatics Association: JAMIA*, 28(7):1591–1599, July 2021.
- [114] Pantelis Linardatos, Vasilis Papastefanopoulos, and Sotiris Kotsiantis. Explainable AI: A Review of Machine Learning Interpretability Methods. *Entropy*, 23(1):18, 2020-12-25.

- [115] Michal Lipschuetz, Joshua Guedalia, Amihai Rottenstreich, Michal Novoselsky Persky, Sarah M. Cohen, Doron Kabiri, Gabriel Levin, Simcha Yagel, Ron Unger, and Yishai Sompolinsky. Prediction of vaginal birth after cesarean deliveries using machine learning. *American Journal of Obstetrics and Gynecology*, 222(6):613.e1–613.e12, 2020.
- [116] Yi Liu. PPGAN Repository. <https://github.com/niklausliu/PPGANs-Privacy-preserving-GANs>.
- [117] Yi Liu, Jialiang Peng, James J.Q. Yu, and Yi Wu. Ppgan: Privacy-preserving generative adversarial network. *Proceedings of the International Conference on Parallel and Distributed Systems - ICPADS*, 2019-Decem(201910212133):985–989, 2019. arXiv: 1910.02007v1 ISBN: 9781728125831.
- [118] S. Lloyd. Least squares quantization in PCM. *IEEE Transactions on Information Theory*, 28(2):129–137, 1982-03.
- [119] Pei-Hsuan Lu, Pang-Chieh Wang, and Chia-Mu Yu. Empirical Evaluation on Synthetic Data Generation with Generative Adversarial Network. In *PROCEEDINGS OF THE 9TH INTERNATIONAL CONFERENCE ON WEB INTELLIGENCE, MINING AND SEMANTICS (WIMS 2019)*, 2019.
- [120] Pei-Hsuan Lu and Chia-Mu Yu. POSTER: A Unified Framework of Differentially Private Synthetic Data Release with Generative Adversarial Network. In *CCS'17: PROCEEDINGS OF THE 2017 ACM SIGSAC CONFERENCE ON COMPUTER AND COMMUNICATIONS SECURITY*, pages 2547–2549. ACM SIGSAC; Assoc Comp Machinery; AT & T Business; Baidu; NSF; CISCO; Internet Finance Authenticat Alliance; Samsung; Univ Texas Dallas; Google; IBM Res; Paloalto Networks; Visa Res; Army Res Off; Nasher Sculpture Ctr, 2017.
- [121] J MacQueen. Classification and analysis of multivariate observations. In *5th Berkeley Symp. Math. Statist. Probability*, pages 281–297. University of California Los Angeles LA USA, 1967.
- [122] Nicole E. Marshall, Rongwei Fu, and Jeanne-Marie Guise. Impact of multiple cesarean deliveries on maternal morbidity: A systematic review. *American Journal of Obstetrics and Gynecology*, 205(3):262.e1–8, 2011.
- [123] F. Martin-Sanchez and K. Verspoor. Big data in medicine is driving big changes. *Yearbook of Medical Informatics*, 9:14–20, August 2014.
- [124] Ewen D. McAlpine, Pamela Michelow, and Turgay Celik. The Utility of Unsupervised Machine Learning in Anatomic Pathology. *American Journal of Clinical Pathology*, 157(1):5–14, 2022-01-06.
- [125] Matthew Michelson and Katja Reuter. The significant cost of systematic reviews and meta-analyses: A call for greater involvement of machine learning to assess the promise of clinical trials. *Contemporary Clinical Trials Communications*, 16:100443, August 2019. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6722281/>.
- [126] Mehdi Mirza and Simon Osindero. Conditional Generative Adversarial Nets, 2014. arXiv: 1411.1784.
- [127] Tom M Mitchell. *Machine learning*, volume 1. McGraw-hill New York, 1997.



- [128] Mahmoud Mohammadi. table-GAN Repository. <https://github.com/mahmoodm2/tableGAN>.
- [129] Raghad Muhiyaddin, Alaa A. Abd-Alrazaq, Mowafa Househ, Tanvir Alam, and Zubair Shah. The Impact of Clinical Decision Support Systems (CDSS) on Physicians: A Scoping Review. In *The Importance of Health Informatics in Public Health during a Pandemic*, pages 470–473. IOS Press, 2020. <https://ebooks.iospress.nl/doi/10.3233/SHTI200597>.
- [130] Emily Muller, Xu Zheng, and Jer Hayes. Evaluation of the Synthetic Electronic Health Records, 2022-10-16.
- [131] Gonzalo Navarro. A guided tour to approximate string matching. *ACM Computing Surveys*, 33(1):31–88, 2001-03.
- [132] Marcel Neunhoeffter. Post-GAN Boosting Repository. <https://github.com/mneunhoe/post-gan-boosting>.
- [133] Marcel Neunhoeffter, Zhiwei Steven Wu, and Cynthia Dwork. Private Post-GAN Boosting, 2020.
- [134] Giovanna Nicora, Francesca Vitali, Arianna Dagliati, Nophar Geifman, and Riccardo Bellazzi. Integrated Multi-Omics Analyses in Oncology: A Review of Machine Learning Methods and Tools. *Frontiers in Oncology*, 10, 2020.
- [135] Beata Nowok, Gillian M. Raab, and Chris Dibben. synthpop: Bespoke creation of synthetic data in R. *Journal of Statistical Software*, 74(11):1–26, 2016.
- [136] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. *34th International Conference on Machine Learning, ICML 2017*, 6:4043–4055, 2017. arXiv: 1610.09585 ISBN: 9781510855144.
- [137] Office for Civil Rights. Guidance Regarding Methods for De-identification of Protected Health Information in Accordance with the Health Insurance Portability and Accountability Act (HIPAA) Privacy Rule. U.S. Department of Health and Human Services, 20 November 2013, 2013.
- [138] Office for Civil Rights. Guidance Regarding Methods for De-identification of Protected Health Information in Accordance with the Health Insurance Portability and Accountability Act (HIPAA) Privacy Rule. U.S. Department of Health and Human Services, 20 November 2013, 2013.
- [139] orsinium. Textdistance: Compute distance between the two texts. <https://github.com/orsinium/textdistance>.
- [140] Venkataraman Palabindala, Amaleswari Pamarthy, and Nageshwar Reddy Jonnalagadda. Adoption of electronic health records and barriers. *Journal of Community Hospital Internal Medicine Perspectives*, 6(5):32643, January 2016. <https://www.tandfonline.com/doi/full/10.3402/jchimp.v6.32643>.
- [141] Trishan Panch, Heather Mattie, and Leo Anthony Celi. The “inconvenient truth” about AI in healthcare. *npj Digital Medicine*, 2(1):77, August 2019. <https://www.nature.com/articles/s41746-019-0155-4>.

- [142] Noseong Park, Mahmoud Mohammadi, Kshitij Gorde, Sushil Jajodia, Hongkyu Park, and Youngmin Kim. Data synthesis based on generative adversarial networks. *Proceedings of the VLDB Endowment*, 11(10):1071–1083, 2018. arXiv: 1806.03384.
- [143] Neha Patki, Roy Wedge, and Kalyan Veeramachaneni. The synthetic data vault. In *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pages 399–410. IEEE, 2016.
- [144] Seyedeh Neelufar Payrovnaziri, Zhaoyi Chen, Pablo Rengifo-Moreno, Tim Miller, Jiang Bian, Jonathan H Chen, Xiuwen Liu, and Zhe He. Explainable artificial intelligence models using real-world electronic health record data: A systematic scoping review. *Journal of the American Medical Informatics Association : JAMIA*, 27(7):1173–1185, 2020-05-17.
- [145] Judea Pearl. Theoretical Impediments to Machine Learning With Seven Sparks from the Causal Revolution.
- [146] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [147] Niels Peek and Pedro Pereira Rodrigues. Three controversies in health data science. *International Journal of Data Science and Analytics*, 6(3):261–269, November 2018. <https://doi.org/10.1007/s41060-018-0109-y>.
- [148] Hang T. T. Phan, Florina Borca, David Cable, James Batchelor, Justin H. Davies, and Sarah Ennis. Automated data cleaning of paediatric anthropometric data from longitudinal electronic health records: Protocol and application to a large patient cohort. *Scientific Reports*, 10(1):10164, June 2020.
- [149] Pordata. Cesarianas nos hospitais (%). [https://www.pordata.pt/Portugal/Cesarianas+nos+hospitais+\(percentagem\)-1985](https://www.pordata.pt/Portugal/Cesarianas+nos+hospitais+(percentagem)-1985).
- [150] POSTER. POSTER Repository. <https://goo.gl/94qyQz>.
- [151] Prayitno, Chi-Ren Shyu, Karisma Trinanda Putra, Hsing-Chung Chen, Yuan-Yu Tsai, K. S. M. Tozammel Hossain, Wei Jiang, and Zon-Yin Shae. A Systematic Review of Federated Learning in the Healthcare Area: From the Perspective of Data Properties and Applications. *Applied Sciences*, 11(23):11191, 2021-01.
- [152] G. Price, M. van Herk, and C. Faivre-Finn. Data mining in oncology: The ukCAT project and the practicalities of working with routine patient data. *Clinical Oncology (Royal College of Radiologists (Great Britain))*, 29(12):814–817, 2017.
- [153] The Synthetic Data Vault Project. CTGAN Repository. <https://github.com/sdv-dev/CTGAN>.
- [154] The Synthetic Data Vault Project. TGAN Repository. <https://github.com/sdv-dev/TGAN>.
- [155] Ross Quinlan. Thyroid Disease. UCI Machine Learning Repository, 1987. DOI: [10.24432/C5D010](https://doi.org/10.24432/C5D010).

- [156] Naresh Sundar Rajan, Ramkiran Gouripeddi, Peter Mo, Randy K. Madsen, and Julio C. Facelli. Towards a content agnostic computable knowledge repository for data quality assessment. *Computer Methods and Programs in Biomedicine*, 177:193–201, August 2019.
- [157] Alvin Rajkomar, Jeffrey Dean, and Isaac Kohane. Machine Learning in Medicine. *New England Journal of Medicine*, 2019-04-03.
- [158] Nimrod Rappoport and Ron Shamir. Multi-omic and multi-view clustering algorithms: Review and cancer benchmark. *Nucleic Acids Research*, 46(20):10546–10562, 2018-11-16.
- [159] Sebastian Raschka. Mlxtend: Providing machine learning and data science utilities and extensions to python’s scientific computing stack. *The Journal of Open Source Software*, 3(24), April 2018.
- [160] Sina Rashidian, Fusheng Wang, Richard Moffitt, Victor Garcia, Anurag Dutt, Wei Chang, Vishwam Pandya, Janos Hajagos, Mary Saltz, and Joel Saltz. SMOOTH-GAN: Towards Sharp and Smooth Synthetic EHR Data Generation. In Martin Michalowski and Robert Moskovitch, editors, *Artificial Intelligence in Medicine*, Lecture Notes in Computer Science, pages 37–48, Cham, 2020. Springer International Publishing.
- [161] Daniele Ravi, Charence Wong, Fani Deligianni, Melissa Berthelot, Javier Andreu-Perez, Benny Lo, and Guang-Zhong Yang. Deep Learning for Health Informatics. *IEEE Journal of Biomedical and Health Informatics*, 21(1):4–21, 2017-01.
- [162] Hanieh Razzaghi, Jane Greenberg, and L. Charles Bailey. Developing a systematic approach to assessing data quality in secondary use of clinical data based on intended use. *Learning Health Systems*, 6(1):e10264, 2022.
- [163] Andrew P. Reimer, Alex Milinovich, and Elizabeth A. Madigan. Data quality assessment framework to assess electronic medical record data for use in research. *International Journal of Medical Informatics*, 90:40–47, June 2016.
- [164] Pedro Pereira Rodrigues, João Araújo, João Gama, and Luís Lopes. A local algorithm to approximate the global clustering of streams generated in ubiquitous sensor networks. *International Journal of Distributed Sensor Networks*, 14(10):155014771880823, 2018-10.
- [165] Seyyed Soroush Rohanizadeh and Mohammad Bameni Moghadam. A Proposed Data Mining Methodology and its Application to Industrial Procedures. *Journal of Industrial Engineering*, 2009.
- [166] Paul R. Rosenbaum. *Causal Inference*. The MIT Press.
- [167] Paul R. Rosenbaum and Donald B. Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983-04-01.
- [168] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019-05.
- [169] H. S. Rugo, V. Diéras, K. A. Gelmon, R. S. Finn, D. J. Slamon, M. Martin, P. Neven, Y. Shparyk, A. Mori, D. R. Lu, H. Bhattacharyya, C. H. U. a. N. G. Bartlett, S. Iyer, S. Johnston,

- J. Ettl, and N. Harbeck. Impact of palbociclib plus letrozole on patient-reported health-related quality of life: Results from the PALOMA-2 trial. *Annals of Oncology: Official Journal of the European Society for Medical Oncology*, 29(4):888–894, April 2018.
- [170] Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach (4th Edition)*. Pearson, 2020.
- [171] D. L. Sackett, W. M. Rosenberg, J. A. Gray, R. B. Haynes, and W. S. Richardson. Evidence based medicine: What it is and what it isn't. 312(7023):71–72.
- [172] Saqib Saleem, Syed Saud Naqvi, Tareq Manzoor, Ahmed Saeed, Naveed ur Rehman, and Jawad Mirza. A strategy for classification of “vaginal vs. cesarean section” delivery: Bi-variate empirical mode decomposition of cardiotocographic recordings. *Frontiers in Physiology*, 10:246, 2019.
- [173] Carsten Oliver Schmidt, Stephan Struckmann, Cornelia Enzenbach, Achim Reineke, Jürgen Stausberg, Stefan Damerow, Marianne Huebner, Borge Schmidt, Willi Sauerbrei, and Adrian Richter. Facilitating harmonized data quality assessments. A data quality framework for observational health research data collections with software implementations in R. *BMC medical research methodology*, 21(1):63, April 2021.
- [174] Sertkaya, Aylin, Birkenbach, Anna, Berlind, Ayesha, and Eyraud, John. Examination of Clinical Trial Costs and Barriers for Drug Development. Technical report, Eastern Research Group, Inc., July 2014. [https://aspe.hhs.gov/sites/default/files/private/pdf/77166/rpt\\_erg.pdf](https://aspe.hhs.gov/sites/default/files/private/pdf/77166/rpt_erg.pdf).
- [175] Jingpu Shi and Beau Norgeot. Learning Causal Effects From Observational Data in Healthcare: A Review and Summary. 9.
- [176] Dennis J. Slamon, Patrick Neven, Stephen Chia, Peter A. Fasching, Michelino De Laurentiis, Seock-Ah Im, Katarina Petrakova, Giulia Val Bianchi, Francisco J. Esteva, Miguel Martín, Arnd Nusch, Gabe S. Sonke, Luis De la Cruz-Merino, J. Thaddeus Beck, Xavier Pivot, Gena Vidam, Yingbo Wang, Karen Rodriguez Lorenc, Michelle Miller, Tetiana Taran, and Guy Jerusalem. Phase III Randomized Study of Ribociclib and Fulvestrant in Hormone Receptor-Positive, Human Epidermal Growth Factor Receptor 2-Negative Advanced Breast Cancer: MONALEESA-3. *Journal of Clinical Oncology: Official Journal of the American Society of Clinical Oncology*, 36(24):2465–2472, August 2018.
- [177] George W. Sledge, Masakazu Toi, Patrick Neven, Joohyuk Sohn, Kenichi Inoue, Xavier Pivot, Olga Burdaeva, Meena Okera, Norikazu Masuda, Peter A. Kaufman, Han Koh, Eva-Maria Grischke, Martin Frenzel, Yong Lin, Susana Barriga, Ian C. Smith, Nawel Bourayou, and Antonio Llombart-Cussac. MONARCH 2: Abemaciclib in Combination With Fulvestrant in Women With HR+/HER2- Advanced Breast Cancer Who Had Progressed While Receiving Endocrine Therapy. *Journal of Clinical Oncology: Official Journal of the American Society of Clinical Oncology*, 35(25):2875–2884, September 2017.
- [178] Theresa Stadler, Bristena Oprisanu, and Carmela Troncoso. Synthetic data – A privacy mirage. *arXiv*, 2020. arXiv: 2011.07018.
- [179] Theresa Stadler, Bristena Oprisanu, and Carmela Troncoso. Synthetic data – A privacy mirage. *arXiv*, 2020.

- [180] Douglas Steinley and Michael J Brusco. Initializing k-means batch clustering: A critical evaluation of several techniques. *Journal of Classification*, 24(1):99–121, 2007.
- [181] Steven Suydam, Bryan A. Liang, Storm Anderson, and Matthew B. Weinger. Patient Safety Data Sharing and Protection From Legal Discovery. *Journal of Medical Regulation*, 93(2):19–25, 2007-06-01.
- [182] Casey Ross Swetlitz, Ike. IBM’s Watson supercomputer recommended ‘unsafe and incorrect’ cancer treatments, internal documents show.
- [183] Uthaipon Tantipongpipat, Chris Waites, Digvijay Boob, Amaresh Ankit Siva, and Rachel Cummings. Differentially Private Synthetic Mixed-Type Data Generation For Unsupervised Learning, 2020. arXiv: cs.LG/1912.03250.
- [184] Achilles Thoma and Felmont F. Eaves, III. A Brief History of Evidence-Based Medicine (EBM) and the Contributions of Dr David Sackett. 35(8):NP261–NP263.
- [185] Ilan E. Timor-Tritsch and Ana Monteagudo. Unforeseen consequences of the increasing rate of cesarean deliveries: Early placenta accreta and cesarean scar pregnancy. A review. *American Journal of Obstetrics and Gynecology*, 207(1):14–29, 2012-07.
- [186] Eric J. Topol. High-performance medicine: The convergence of human and artificial intelligence. 25(1):44–56.
- [187] Amirsina Torfi. corGAN Repository. <https://github.com/astorfi/cor-gan>.
- [188] Amirsina Torfi. DRP-CGAN Repository.
- [189] Amirsina Torfi and Edward A. Fox. CorGAN: Correlation-capturing convolutional generative adversarial networks for generating synthetic healthcare records. *arXiv*, 2020. arXiv: 2001.09346.
- [190] Amirsina Torfi, Edward A. Fox, and Chandan K. Reddy. Differentially Private Synthetic Medical Data Generation using Convolutional GANs. *arXiv:2012.11774 [cs]*, December 2020. arXiv: 2012.11774; <https://web.archive.org/web/20210618105126/https://arxiv.org/abs/2012.11774>.
- [191] Debu Tripathy, Seock-Ah Im, Marco Colleoni, Fabio Franke, Aditya Bardia, Nadia Harbeck, Sara A. Hurvitz, Louis Chow, Joohyuk Sohn, Keun Seok Lee, Saul Campos-Gomez, Rafael Villanueva Vazquez, Kyung Hae Jung, K. Govind Babu, Paul Wheatley-Price, Michelino De Laurentiis, Young-Hyuck Im, Sherko Kuemmel, Nagi El-Saghir, Mei-Ching Liu, Gary Carlson, Gareth Hughes, Ivan Diaz-Padilla, Caroline Germa, Samit Hirawat, and Yen-Shen Lu. Ribociclib plus endocrine therapy for premenopausal women with hormone-receptor-positive, advanced breast cancer (MONALEESA-7): A randomised phase 3 trial. *The Lancet. Oncology*, 19(7):904–915, July 2018.
- [192] Anup Tuladhar, Sascha Gill, Zahinoor Ismail, and Nils D. Forkert. Building machine learning models without sharing patient data: A simulation-based analysis of distributed learning by ensembling. *Journal of Biomedical Informatics*, 106:103424, 2020.
- [193] Zahid Ullah, Farrukh Saleem, Mona Jamjoom, and Bahjat Fakieh. Reliable prediction models based on enriched data for identifying the mode of childbirth by using machine learning methods: Development study. *Journal of Medical Internet Research*, 23(6):e28856, 2021.

- [194] Belen Vega-Marquez, Cristina Rubio-Escudero, Jose C Riquelme, and Isabel Nepomuceno-Chamorro. Creation of Synthetic Data with Conditional Generative Adversarial Networks. In *14TH INTERNATIONAL CONFERENCE ON SOFT COMPUTING MODELS IN INDUSTRIAL AND ENVIRONMENTAL APPLICATIONS (SOCO 2019)*, volume 950, pages 231–240. Startup Ole; IEEE SMC Spanish Chapter, 2020.
- [195] Robert A. Verheij, Vasa Curcin, Brendan C. Delaney, and Mark M. McGilchrist. Possible Sources of Bias in Primary Care Electronic Health Record Data Use and Reuse. *Journal of Medical Internet Research*, 20(5):e9134, May 2018.
- [196] Sunil Verma, Cynthia Huang Bartlett, Patrick Schnell, Angela M. DeMichele, Sherene Loi, Jungsil Ro, Marco Colleoni, Hiroji Iwata, Nadia Harbeck, Massimo Cristofanilli, Ke Zhang, Alexandra Thiele, Nicholas C. Turner, and Hope S. Rugo. Palbociclib in Combination With Fulvestrant in Women With Hormone Receptor-Positive/HER2-Negative Advanced Metastatic Breast Cancer: Detailed Safety Analysis From a Multicenter, Randomized, Placebo-Controlled, Phase III Study (PALOMA-3). *The Oncologist*, 21(10):1165–1175, October 2016.
- [197] Sebastiano Vigna. A Weighted Correlation Index for Rankings with Ties. In *Proceedings of the 24th International Conference on World Wide Web, WWW '15*, pages 1166–1176. International World Wide Web Conferences Steering Committee, 2015-05-18.
- [198] Angela G. Villanueva, Robert Cook-Deegan, Barbara A. Koenig, Patricia A. Deverka, Erika Versalovic, Amy L. McGuire, and Mary A. Majumder. Characterizing the Biomedical Data-Sharing Landscape. *The Journal of law, medicine & ethics : a journal of the American Society of Law, Medicine & Ethics*, 47(1):21–30, 2019-03.
- [199] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020.
- [200] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C. J. Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, and Paul van Mulbregt. SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nature Methods*, 17(3):261–272, 2020-03.
- [201] Manhar Walia, Brendan Tierney, and Susan McKeever. Synthesising tabular data using wasserstein conditional gans with gradient penalty (wsgan-gp). In *Irish Conference on Artificial Intelligence and Cognitive Science*, 2020.
- [202] Muhammad F. Walji. Electronic Health Records and Data Quality. *Journal of Dental Education*, 83(3):263–264, March 2019.

- [203] Abigail Walker and Pavol Surda. Unsupervised Learning Techniques for the Investigation of Chronic Rhinosinusitis. *The Annals of Otology, Rhinology, and Laryngology*, 128(12):1170–1176, 2019-12.
- [204] Shannon C. Walker, Benjamin French, Ryan P. Moore, Henry J. Domenico, Jonathan P. Wanderer, Amanda S. Mixon, C. Buddy Creech, Daniel W. Byrne, and Allison P. Wheeler. Model-Guided Decision-Making for Thromboprophylaxis and Hospital-Acquired Thromboembolic Events Among Hospitalized Children and Adolescents: The CLOT Randomized Clinical Trial. 6(10):e2337789.
- [205] Fei Wang and Anita Preininger. AI in Health: State of the Art, Challenges, and Future Directions. *Yearbook of Medical Informatics*, 28(1):16–26, 2019-08.
- [206] Yining Wang, Liwei Wang, Yuanzhi Li, Di He, Tie-Yan Liu, and Wei Chen. A theoretical analysis of ndcg type ranking measures, 2013.
- [207] Stefanie Warnat-Herresthal, Hartmut Schultze, Krishnaprasad Lingadahalli Shastry, Sathyanarayanan Manamohan, Saikat Mukherjee, Vishesh Garg, Ravi Sarveswara, Kristian Händler, Peter Pickkers, N. Ahmad Aziz, Sofia Ktena, Florian Tran, Michael Bitzer, Stephan Ossowski, Nicolas Casadei, Christian Herr, Daniel Petersheim, Uta Behrends, Fabian Kern, Tobias Fehlmann, Philipp Schommers, Clara Lehmann, Max Augustin, Jan Rybníček, Janine Altmüller, Neha Mishra, Joana P. Bernardes, Benjamin Krämer, Lorenzo Bonaguro, Jonas Schulte-Schrepping, Elena De Domenico, Christian Siever, Michael Kraut, Milind Desai, Bruno Monnet, Maria Saridaki, Charles Martin Siegel, Anna Drews, Melanie Nuesch-Germano, Heidi Theis, Jan Heyckendorf, Stefan Schreiber, Sarah Kim-Hellmuth, Jacob Nattermann, Dirk Skowasch, Ingo Kurth, Andreas Keller, Robert Bals, Peter Nürnberg, Olaf Rieß, Philip Rosenstiel, Mihai G. Netea, Fabian Theis, Sach Mukherjee, Michael Backes, Anna C. Aschenbrenner, Thomas Ulas, Monique M. B. Breteler, Evangelos J. Giamarellos-Bourboulis, Matthijs Kox, Matthias Becker, Sorin Cheran, Michael S. Woodacre, Eng Lim Goh, and Joachim L. Schultze. Swarm Learning for decentralized and confidential clinical machine learning. *Nature*, 594(7862):265–270, June 2021.
- [208] William Webber, Alistair Moffat, and Justin Zobel. A similarity measure for indefinite rankings. *ACM Transactions on Information Systems*, 28(4):20:1–20:38, 2010-11-23.
- [209] Nicole Gray Weiskopf and Chunhua Weng. Methods and dimensions of electronic health record data quality assessment: Enabling reuse for clinical research. 20(1):144–151.
- [210] Nicole Gray Weiskopf and Chunhua Weng. Methods and dimensions of electronic health record data quality assessment: Enabling reuse for clinical research. *Journal of the American Medical Informatics Association*, 20(1):144–151, January 2013.
- [211] Chunhua Weng. Clinical data quality: A data life cycle perspective. *Biostatistics & Epidemiology*, 4(1):6–14, January 2020.
- [212] William Wolberg, Olvi Mangasarian, Nick Street, and W. Street. Breast Cancer Wisconsin (Diagnostic). UCI Machine Learning Repository, 1995. DOI: [10.24432/C5DW2B](https://doi.org/10.24432/C5DW2B).
- [213] World Health Organization Human Reproduction Programme, 10 April 2015. WHO statement on caesarean section rates. *Reproductive Health Matters*, 23(45):149–150, 2015.
- [214] Liyang Xie, Kaixiang Lin, Shu Wang, Fei Wang, and Jiayu Zhou. Differentially private generative adversarial network. *arXiv*, 2018. arXiv: 1802.06739 ISBN: 1234567245.



- [215] Jie Xu, Benjamin S. Glicksberg, Chang Su, Peter Walker, Jiang Bian, and Fei Wang. Federated Learning for Healthcare Informatics. *Journal of Healthcare Informatics Research*, pages 1–19, 2020-11-12.
- [216] Jie Xu, Benjamin S. Glicksberg, Chang Su, Peter Walker, Jiang Bian, and Fei Wang. Federated Learning for Healthcare Informatics. *Journal of Healthcare Informatics Research*, 5(1):1–19, 2021.
- [217] Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Modeling tabular data using conditional GAN. *arXiv*, 32(NeurIPS), 2019. arXiv: 1907.00503.
- [218] Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Modeling Tabular data using Conditional GAN. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [219] Lei Xu and Kalyan Veeramachaneni. Synthesizing tabular data using generative adversarial networks. *arXiv*, November 2018. arXiv: 1811.11264.
- [220] Qiantong Xu, Gao Huang, Yang Yuan, Chuan Guo, Yu Sun, Felix Wu, and Kilian Weinberger. An empirical study on evaluation metrics of generative adversarial networks. *arXiv:1806.07755 [cs, stat]*, August 2018. arXiv: 1806.07755; <https://web.archive.org/web/20200604163128/https://arxiv.org/abs/1806.07755>.
- [221] Yan Cheng Yang, Saad Ul Islam, Asra Noor, Sadia Khan, Waseem Afsar, and Shah Nazir. Influential Usage of Big Data and Artificial Intelligence in Healthcare. *Computational and Mathematical Methods in Medicine*, 2021:5812499, 2021-09-06.
- [222] Jinsung Yoon, Lydia N Drumright, and Mihaela van der Schaar. Anonymization Through Data Synthesis Using Generative Adversarial Networks (ADS-GAN). *IEEE JOURNAL OF BIOMEDICAL AND HEALTH INFORMATICS*, 24(8):2378–2388, 2020.
- [223] Yili Zhang and Güneş Koru. Understanding and detecting defects in healthcare administration data: Toward higher data quality to better support healthcare operations and decisions. *Journal of the American Medical Informatics Association: JAMIA*, 27(3):386–395, March 2020.
- [224] M. Zwitter and M. Soklic. Primary Tumor. UCI Machine Learning Repository, 1988. DOI: [10.24432/C5WK5Q](https://doi.org/10.24432/C5WK5Q).