



The Ottawa | L'Hôpital Hospital d'Ottawa



#MEDINF023

Randi E. Foraker, Ph.D., M.A. Adam Wilcox, Ph.D. Alan Forster, M.D., M.Sc. Zachary Abrams, Ph.D. Jon D. Morrow, M.D., M.A., M.B.A.

Synthetic Data

Expanding Information Accessibility:



Session Objectives

- The panellists will explore the current and future state of synthetic data, highlighting their ability to support data sharing, to address privacy and confidentiality, and to advance national and international initiatives.
- Panellists will address the following topics:
 - Statistical validation
 - Privacy validation
 - Enhancing data sharing and collaboration
 - Industry partnerships
 - Use in education and training programs
 - Methodology and process of synthetic data generation



Panellists



Zachary Abrams, Ph.D. Washington University in St. Louis



Randi E. Foraker, Ph.D. Washington University in St. Louis







Jon D. Morrow, M.D. New York University



Adam Wilcox, Ph.D. Washington University in St. Louis



Session Agenda

	Dr Morrow	" Synthetic Data Primer" Synthetic data background, definitions, methods
	Prof Foraker	"Validation of Synthetic Data" Statistical validation, privacy, and research use-cases
	Dr Forster	"Creating a Learning System" Democratisation of data; collaborative partnerships
	Prof Abrams Prof Foraker Dr Forster Dr Morrow Prof Wilcox	Panel discussion D&A
aTheln	stituteDH #MEDINE073	





Synthetic Data Primer



Jon D. Morrow, MD, MA, MBA, FACOG

Clinical Associate Professor Department of Obstetrics & Gynecology New York University School of Medicine



@TheInstituteDH #MEDINF023





Disclosure

Jon D. Morrow, M.D., was previously Senior Vice President at MDClone, Inc., a health IT software company whose synthetic data product was used in some of the referenced studies.

Why use synthetic health data?

- Personal health data are highly protected and sensitive.
- Health data contain a vast amount of valuable information.
- Data-to-knowledge work is often done by people not privy to protected information (e.g., non-clinicians, consortium members, and external commercial partners).



What are synthetic data?

- Maintain the utility, statistical properties, correlations, and higher-order relationships of real data sets without containing or exposing the members of the original set.
- Same format and suitability for analysis as the original.
- No one-to-one correspondence between members of the original and synthetic data sets.
- A form of data anonymisation, but:

Synthetic data ≠ Deidentified data









@TheInstituteDH

#MEDINFUZ3

DRIGINAL





Ν	2,255
Age, mean ± SD	32.5 ± 6.7 years
Height, mean ± SD	171.1 ± 6.35 cm
Gender	62.1% male 36.9% female 1.0% non-binary
DM	8.2%
BMI ≥ 30 kg/m²	35.1%
DM among BMI < 30 DM among BMI ≥ 30	6.9% 12.0%
M	0.057
Ν	2,257
N Age, mean ± SD	2,257 32.4 ± 6.7 years
N Age, mean ± SD Height, mean ± SD	2,257 32.4 ± 6.7 years 170.9 ± 6.32 cm
N Age, mean ± SD Height, mean ± SD Gender	2,257 32.4 ± 6.7 years 170.9 ± 6.32 cm 62.0% male 37.1% female 0.9% non-binary
N Age, mean ± SD Height, mean ± SD Gender DM	2,257 32.4 ± 6.7 years 170.9 ± 6.32 cm 62.0% male 37.1% female 0.9% non-binary 8.2%
N Age, mean ± SD Height, mean ± SD Gender DM BMI ≥ 30 kg/m ²	2,257 32.4 ± 6.7 years 170.9 ± 6.32 cm 62.0% male 37.1% female 0.9% non-binary 8.2% 35.3%

SD = Standard deviation; DM = Diabetes mellitus; BMI = Body mass index.



@TheInstituteDH

#MEDINF023

DRIGINAL





Age, mean ± SD 32.5 ± 6.7 years Height, mean ± SD 171.1 ± 6.35 cm 62.1% male Gender 36.9% female 1.0% non-binary DM $BMI \ge 30 \text{ kg/m}^2$ 35.1% DM among BMI < 30 DM among BMI ≥ 30 Ν 2,255 30.9 ± 5.8 years Age, mean ± SD Height, mean ± SD 171.1 ± 6.35 cm 62.1% male Gender 36.9% female 1.0% non-binary DM 8.2% $BMI \ge 30 \text{ kg/m}^2$ 35.1% DM among BMI < 30 6.9%

D = Standard deviation; DM = Diabetes mellitus; BMI = Body mass index

12.0%

DM among BMI \ge 30



@TheInstituteDH

#MEDINF023

Simulation vs computation

- <u>Simulated synthetic data</u>: Probabilistic synthesis to create large data sets, useful for simulation, systems testing, training, and other uses *(e.g., Synthea*?).
- <u>Computationally derived synthetic data</u>: Novel data set, usually (but not necessarily) approximately the same size as the original, populated with new data points to match the original's statistical properties (e.g., MDClone ADAMSTM).



Synthetic data for Al



- ML models trained on computationally derived synthetic data are equally valid as models trained on the original source data from which the synthetic data were derived.
- This allows ML models to be trained without exposing the original patients' personal information and without compromising patient privacy.





Validation of Synthetic Data



Randi Foraker, PhD, MA, FAHA, FAMIA, FACMI

Director, Center for Population Health Institute for Informatics (12) Washington University in St. Louis



@TheInstituteDH #MEDINF023



@TheInstituteDH #MEDINF023

C

A.

Validation studies: Questions answered

Privacy Studies (n=5)	Statistical Studies (n=6)
Does anyone look the same?	Does it look the same?
Can you automatically identify people?	Does it work the same?
Can you manually identify people?	Does it work better than other approaches?
It is a problem if people are identified?	What can we do differently with synthetic data?
Does it solve the problem of linked data sets?	

@TheInstituteDH #MEDINF023



Research and Applications

Spot the difference: comparing results of analyses from real patient data and synthetic derivatives

Randi E. Foraker (*),^{1,2} Sean C. Yu,² Aditi Gupta,² Andrew P. Michelson,³ Jose A. Pineda Soto,⁴ Ryan Colvin,^{2,4} Francis Loh,⁵ Marin H. Kollef,³ Thomas Maddox,⁶ Bradley Evanoff,¹ Hovav Dror,⁷ Noa Zamstein,⁷ Albert M. Lai (*),^{1,2} and Philip R.O. Payne (*)^{1,2}

"Division of General Medical Sciences, Department of Medicine, School of Medicine, Washington University in St. Losis, St. Losis, MSaouri, USA, "Department of Medicine, Institute for Informatics, School of Medicine, Nathington University in St. Losis, St. Losis, MSaouri, USA, "Division of Pulmonary and Critical Care Medicine, Department of Medicine, School of Medicine, Weshington University in St. Losis, St. Losis, MSaouri, USA, "Division of Critical Care Medicine, Department of Ametical-Restrictions, St. School of Medicine, Washington University in St. Losis, St. Losis, MSaouri, USA, "Division of Critical Care Medicine, Department of Ametical-Restriction, Children's Hospital of Los Angeles, Carlinonia, USA, "School of Medicine, Washington University in St. Losis, St. Losis, Msaouri, USA, "Metalthcare Innovation Lab, BUC Healthcare, School of Medicine, Washington University in St. Losis, S. Losis, Msaouri, USA, "Metalthcare Innovation Lab, BUC Healthcare, School of Medicine, Washington University in St. Losis, St. Losis, Msaouri, USA, "Metalthcare Innovation Lab, BUC Healthcare, School of Medicine, Washington University in St. Losis, St. Losis, Msaouri, USA, "Metalthcare Innovation Lab, BUC Healthcare, School of Medicine, Washington University in St. Losis, St. Losis, Msaouri, USA, "Metalthcare Innovation Lab, BUC Healthcare, School of Medicine, Washington University in St. Losis, St. Losis, Msaouri, USA, Martine, St. Restriction, Stal, "School", Stal, Chabit, Msaouri, USA, Matthicare, School of Medicine, Washington University in St. Losis, St. Losis, Msaouri, USA, Martine, St. Restriction, Stal, "School", Stal, School, Stal, Stal, St. Losis, Msaouri, USA, Matthicare, School of Medicine, Washington University in St. Losis, St. Losis, Msaouri, USA, Matthicare, School I, Bert Sheve, Israel

Corresponding Author: Randi Forakar, PhD, MA, FAHA, FAMIA, Associate Professor, Institute for Informatics (I2), Director, Center for Population Health Informatics a I2, Washington University in St. Louis, School of Medicine, 600 S. Taylor Avenue, Suite 102, Campus Box 8102, St. Louis, MO 63110, USA (http://informatics.wustledu)

Received 14 August 2020; Revised 14 October 2020; Accepted 20 October 2020

ABSTRACT

Background: Synthetic data may provide a solution to researchers who wish to generate and share data in support of precision healthcare. Recent advances in data synthesis enable the creation and analysis of synthetic derivatives as if they were the original data; this process has significant advantages over data deidentification. Objectives: To assess a big-data platform with data-synthesizing capabilities (MDClone Ltd., Beer Sheva, Israel) for its ability to produce data that can be used for research purposes while obviating privacy and confidentiality concerns.

Methods: We explored three use cases and tested the robustness of synthetic data by comparing the results of analyses using synthetic derivatives to analyses using the original data using traditional statistics, machine learning approaches, and spatial representations of the data. We designed these use cases with the purpose of conducting analyses at the observation level (Use Case 1), patient cohorts (Use Case 2), and population-level data (Use Case 3).

Results: For each use case, the results of the analyses were sufficiently statistically similar (P>0.05) between the synthetic derivative and the real data to draw the same conclusions.

Discussion and conclusion: This article presents the results of each use case and outlines key considerations for the use of synthetic data, examining their role in clinical research for faster insights and improved data sharing in support of precision healthcare.

Key words: synthetic data, protected health information, precision health care, electronic health records and systems, data analysis

© The Author(s) 2029. Published by Oxford University Press on behalf of the American Medical Informatics Association. This is an Open Access ancide distributed under the terms of the Creative Commons Attribution License Intp://creativecommons.org/licenses/by/4.0/), which permits unstrictical reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Performance metrics of sepsis prediction models

Training set		Real		Synthetic	Synthetic	
Testing set		Real		Synthetic	Real	
Train	Accuracy	0.845	Π	0.869	0.852	
	Precision	0.803		0.840	0.812	
	Recall	0.704		0.758	0.719	
	FI	0.750		0.797	0.763	
	AUROC	0.809		0.842	0.818	
5-fold cross-	Accuracy	0.795		0.802	0.799	
validation	Precision	0.712		0.73	0.723	
	Recall	0.637		0.67	0.639	
	FI	0.672		0.69	0.678	
	AUROC	0.855		0.86	0.847	
Test	Accuracy	0.834		0.833	0.834	
	Precision	0.811		0.759	0.829	
	Recall	0.677		0.678	0.654	
	FI	0.738		0.716	0.731	
	AUROC	0.887		0.885	0.892	

JAMIA Open, Volume 3, Issue 4, December 2020, Pages 557–566, doi:10.1093/jamiaopen/ooaa060

@TheInstituteDH



Research and Applications

Spot the difference: comparing results of analyses from real patient data and synthetic derivatives

Randi E. Foraker (*), ^{1,2} Sean C. Yu,² Aditi Gupta,² Andrew P. Michelson,³ Jose A. Pineda Soto,⁴ Ryan Colvin,^{2,4} Francis Loh,⁵ Marin H. Kollef,³ Thomas Maddox,⁶ Bradley Evanoff, ¹ Hovav Dror,⁷ Noa Zamstein,⁷ Albert M. Lai (*),^{1,2} and Philip R.O. Payne (*).^{1,2}

¹Division of General Medical Sciences, Department of Medicine, School of Medicine, Washington University in St. Louis, St. Louis, Masouri, USA, ¹Department of Medicine, Institute for Informatics, School of Medicine, Nathington University in St. Louis, St. Louis, Masouri, USA, ¹Division of Pulmonary and Critical Care Medicine, Department of Medicine, School of Medicine, Weshington University in St. Louis, St. Louis, Masouri, USA, ¹Division of Critical Care Medicine, Department of Medicine, School of Medicine, Washington University in St. Louis, St. Louis, Masouri, USA, ¹Division of Critical Care Medicine, Department of Anesthesiology and Critical Care Medicine, USA, ¹School of Medicine, Washington University in St. Louis, St. Louis, Masouri, USA, ¹Mealthcare Innovation Lab, BUC Healthcare, School of Medicine, Washington University in St. Louis, St. Louis, Masouri, USA, ¹Mealthcare Innovation Lab, BUC Healthcare, School of Medicine, Washington University in St. Louis, St. Louis, Masouri, USA, ¹Mealthcare Lineovation Lab, BUC Healthcare, School of Medicine, Washington University in St. Louis, St. Louis, Masouri, USA, ¹Mealthcare Envoyation Lab, BUC Healthcare, School of Medicine, Washington University in St. Louis, St. Louis, Masouri, USA, ¹Mealthcare Envoyation Lab, BUC Healthcare, School of Medicine, Washington University in St. Louis, St. Louis, Masouri, USA, ¹Mealthcare Envoyation Lab, BUC Healthcare, School of Medicine, Washington University in St. Louis, St. Louis, Masouri, USA and ¹MDCDorea.

Corresponding Author: Randi Foraker, PhD, MA, FAHA, FAMIA, Associate Professor, Institute for Informatics (12), Director, Center for Population Health Informatics at 12, Washington University in St. Louis, School of Medicine, 600 S. Taylor Avenue, Suite 102, Campus Box 8102, St. Louis, MG 63110, USA (http://nformatics.wustledu)

Received 14 August 2020; Revised 14 October 2020; Accepted 20 October 2020

ABSTRACT

Background: Synthetic data may provide a solution to researchers who wish to generate and share data in support of precision healthcare. Recent advances in data synthesis enable the creation and analysis of synthetic derivatives as if they were the original data; this process has significant advantages over data diedinetification. Objectives: To assess a big-data platform with data-synthesizing capabilities (MDClone Ltd., Beer Sheva, Israel) for its ability to produce data that can be used for research purposes while obviating privacy and confidentiality concerns.

Methods: We explored three use cases and tested the robustness of synthetic data by comparing the results of analyses using synthetic derivatives to analyses using the original data using traditional statistics, machine learning approaches, and spatial representations of the data. We designed these use cases with the purpose of conducting analyses at the observation level (Use Case 1), patient cohorts (Use Case 2), and population-level data (Use Case 3).

Results: For each use case, the results of the analyses were sufficiently statistically similar (P>0.05) between the synthetic derivative and the real data to draw the same conclusions.

Discussion and conclusion: This article presents the results of each use case and outlines key considerations for the use of synthetic data, examining their role in clinical research for faster insights and improved data sharing in support of precision healthcare.

Key words: synthetic data, protected health information, precision health care, electronic health records and systems, data analysis

© The Author(s) 2020. Published by Odrod University Press on behalf of the American Medical Informatics Association.
This is an Open Access article distributed under the terms of the Creative Commons Art/buton License (http://creativecommons.org/licenses/by/4.0/), which permits
unversited reverse, distribution, and reproduction in any mediation, provided the original work is properly cited.
1

Chlamydia rates (per 100 000 persons) by zip code: real (left) versus synthetic (right) data.



*Darker color indicates a higher rate

JAMIA Open, Volume 3, Issue 4, December 2020, Pages 557–566, doi:10.1093/jamiaopen/ooaa060

Prediction performance for the two models by receiver operating characteristic curves (A, C) and precision-recall curves (B, D) by using original and synthetic data.



J Med Internet Res. 2021;23(10):e30697. doi:10.2196/30697. PMID: 34559671; PMCID: PMC8491642.

JOURNAL OF MEDICAL INTERNET RESEARCH

Foraker et al

Original Paper

The National COVID Cohort Collaborative: Analyses of Original and Computationally Derived Electronic Health Record Data

Randi Foraker¹², MA, PhD: Aixia Guo³, PhD; Jason Thomas⁴, BS; Noa Zamstein⁴, MSe, PhD; Philip RO Payne¹², PhD: Adam Wilcox³, PhD: N3C Collaborative³

"Division of casars' Medical Sectors, School of Medicine, Washington Divisionity in School, School of Medical States Institute for Information, School of Medicine, Washington Divisionity in School St. Louis, MO, United States "Degrement of University" (School of Medicine, Liviewsky of Washington: Seculity WA, United States "MitChion, Lin, Bass Secure, 1971 "See Advancements"

Corresponding Author:

Randi Fornker, MA, PhD Division of General Madiaal Sciences School of Moffleme Washington Linkershy in St. Louis 600 S. Luyler, Avenue, Suite 102 Campar Box 85, 100 Linited States Pronet: 1311/25211 Fast: 134/273/1390 Fast: 134/273/1390

Abstract

Background: Compatationally derived ("synthetic") data can enable the creation and analysis of clinical, laboratory, and diagnostic data as if they were the original determine health record data. Synthetic data can support data sharing to answer critical research questions to address the COVID 19 pandemic.

Objective: We aim to compare the results from analyses of synthetic data to those from original data and assess the strengths and limitations of loveraging computationally derived data for research purposes.

Methods: We usual the Valiend COVID Coher Collaborative's insures of MICCone, the gluto platform with duta-synthesizing capabilities (MDCiane Lida: We downloaded electronic health neceed data from 34 National COVID Coher Collaborative minimum program in totacil three use cross, including (1) copluring the dimbutions of log features of the COVID-9 positive coher: C) utaining and testing predictive models for seasoling the file of administra aroung these patients; and (5) determining coportal and removed COVID 19-electronic means and universes, and contractiving the equivalence in envice. We compared the results from synthesic data to those from original data using traditional statistics, mechane learning approaches, and temporal and spatial representations of the data.

Regults: For each ore case, the results of the synthetic data analyses recervisitly mimicked bases of the original data wash that the distributions of the data wave similar and the predictive models demonstrated comparable performance. Although the synthetic and original data wave similar and the predictive models the non-service scenarios that included an odds ratio on either vide of the mill in introtraviable analyses (037 vs. 101) and differences in the magnitude of opidemic current constructed for zip codes with like peptidinion currents.

Conclusions: This paper presents the results of each use case and outlines key considerations for the use of synthetic data, examining their role in collaborative research for faster insights.

@TheInstituteDH #MEDINF023

JOURNAL OF MEDICAL INTERNET RESEARCH

Foraker et al

Original Paper

The National COVID Cohort Collaborative: Analyses of Original and Computationally Derived Electronic Health Record Data

Randi Foraker^{1,2}, MA. PhD: Aixia Gue², PhD; Jason Thomas⁶, BS; Noa Zamstein¹, MSe, PhD; Philip RO Payne^{1,2}, PhD; Adam Wilcox³, PhD; N3C Collaborative³

¹Dirician of Gustari Medical Sciences, School of Medicine, Wackingon Enterrity in Sci. Louis, 8:1. Louis, 8:0, United States Instanto for informatics. School of Medicine, Wachington University in Sci. Louis, 3:1. Louis, 3:00. United States "Department of Universities" and Water Informatic School of Medicine, University of Wachington: Scattle Web, United States "MEGRON Edu, Fase Scass, Ford

Corresponding Author:

Randi Forenker, MA, PhD Division of Cherneri Medical Sciences School of Medicine Wachington Thivensity in St. Louis 600 S. Lycler Avenue, Sulte 102 Campate Jox 8102 Linited States Promet: 1311 775 2311 Fast: 134 273 1390 Fast: 134 273 1390

Abstract

Background: Computationally derived ("synthetic") data can enable the creation and analysis of clinical, laboratory, and diagnosic data us if they were the original electronic health record data. Synthetic data can support data sharing to snever critical research questions to address the COVID 19 pandenile.

Objective: We aim to compare the results from analyses of synthetic data to those from original data and assess the strengths and limitations of laveraging computationally derived data for research purposes.

Methods: We used the Valleed QOVID Coher Collaborative's instance of MICCase, all glue platform with indus-synthesizing capabilities MIDCiane. Eds. We downloaded elexinish enablin recend data from 64 Autoard COMID Coher Collaborative midiational partners and tooled these accurses, including (1) exploring the dismithations of key features of the COVID-19-positive coher. Eds. Juraning and testing predictive models for seassing the file of administion among these patients; and (5) determining agoptabilit and resulting predictive models for seassing the file of administion among these patients; and (5) determining agoptabil and resulting predictive models for seassing the file of administration gives relations; and reprint line results from synthetic data to those from original data using traditional statistics, machine learning approaches, and remporal and spatial representitions of the data.

Regults: For each use case, the results of the synthetic data manyces recervfully mimited bases of the original data wash that the data bits of the data wave similar and the predictive models demonstrated comparable performance. Although the synthesis and original data wave similar and the predictive models the non-service scenarios that included an odds ratio on either side of the mill in untravalable analyses (037 vs. 1.01) and differences in the magnitude of opdemic curves constructed for zip codes with here population curves.

Conclusions: This paper presents the results of each use case and outlines key considerations for the use of synthetic data, examining their role in collaborative research for faster insights.

All sites positive tests (cases)



J Med Internet Res. 2021;23(10):e30697. doi:10.2196/30697. PMID: 34559671; PMCID: PMC8491642.

TheInstituteDH #MEDINF023

"...empowers researchers to produce valid results, over a short period of time, while protecting patient privacy."





Creating a Learning System

The Ottawa | L'Hôpital Hospital d'Ottawa

Alan Forster, MD, MSc

Executive Vice President of Innovation and Quality

The Ottawa Hospital



Evidence based decision making in healthcare: what is the problem?

Technical

- Privacy
- Complexity
- Quality
- Completeness



Cultural

- Decision making
- Priorities
- Action tracking
- Trust





@TheInstituteDH #MEDINF023

Data Democratization

Human Data:





A new way of working together





Collaboration at scale







8 – 12 JULY 2023 | SYDNEY, AUSTRALIA

Panel Discussion



Moderator: Randi E. Foraker, Ph.D. Washington University in St. Louis



Alan Forster, M.D. The Ottawa Hospital



Jon D. Morrow, M.D. New York University



Zachary Abrams, Ph.D. Washington University in St. Louis



Adam Wilcox, Ph.D. Washington University in St. Louis



8 – 12 JULY 2023 | SYDNEY, AUSTRALIA









randi.foraker@wustl.edu



@TheInstituteDH #MEDINF023

RSEWES DS



Formal definition of synthetic data

- "Synthetic data are microdata records created to **improve data utility** while **preventing disclosure** of confidential respondent information."
- "Synthetic data is created by statistically modeling original data and then using those models to generate new data values that reproduce the original data's statistical properties."

Philpott D. A Guide to Federal Terms and Acronyms. Bernan Press, 2017:184.



DRIGINAL





N	2,255
Age, mean ± SD	32.5 ± 6.7 years
Height, mean ± SD	171.1 ± 6.35 cm
Gender	62.1% male 36.9% female 1.0% non-binary
DM	8.2%
$BMI \ge 30 \text{ kg/m}^2$	35.1%
DM among BMI < 30 DM among BMI ≥ 30	6.9% 12.0%
N	2,255

	2,200
Age, mean ± SD	30.9 ± 5.8 years
Height, mean ± SD	171.1 ± 6.35 cm
Gender	62.1% male 36.9% female 1.0% non-binary
DM	8.2%
BMI ≥ 30 kg/m²	35.1%
DM among BMI < 30 DM among BMI ≥ 30	6.9% 12.0%

SD = Standard deviation; DM = Diabetes mellitus; BMI = Body mass index.



#MEDINF023



@TheInstituteDH #MEDINF023

Maps: NYC Dept. of Health and Mental Hygiene. NYC STI Surveillance Report, 2021.

Computational derivation overview

- First, derive statistical models from the original data set.
- Then, sample novel synthetic data points to fit the models.





Synthetic Patient Data

Figure adapted from: Foraker RE, Yu SC, Gupta A, et al. Spot the difference: comparing results of analyses from real patient data and synthetic derivatives. *JAMIA Open* 2020;3:557-66.



Computational derivation methodology

- For <u>categorical variables</u>, to mitigate the potential for an inference attack due to the finite number of categories:
 - Group synthetic individuals who share identical categories.
 - If any group contains <ĸ members, censor some discrete values until all groups contain ≥ĸ members.
 - Create clusters of < κ individuals, minimizing the scaled Euclidean distance between data points.
 - Replace each cluster's numeric variables with an alternate matrix with similar statistical properties, preserving statistical characteristics for every pair of variables within each cluster.
 - The alternate matrix is selected randomly from the unlimited number of alternatives, resulting in an irreversible transformation.
- To protect against a difference attack, slightly alter the population size.

Thomas JA, Foraker RE, Zamstein N, Morrow JD, Payne PRO, Wilcox A. Demonstrating an approach for evaluating synthetic geospatial and temporal epidemiologic data utility: results from analyzing >1.8 million SARS-CoV-2 tests in the United States National COVID Cohort Collaborative (N3). *J Amer Med inform Assoc* 2022;29:1350-65.



RSEW SIDS Rend

https://informatics.wustl.edu/mdclone/

Lessons learned

- Synthetic data benefits:
 - Enhancing AI/ML training opportunities
 - Accelerating research
 - Facilitating data sharing
 - Expanding data access to community partners / health dashboards
- Next steps:
 - Expanding upon 'information gain' analyses
 - Establishing an STL data hub



The Use of Synthetic Electronic Health Record Data and Deep Learning to Improve Timing of High-Risk Heart Failure Surgical Intervention by Predicting Proximity to Catastrophic Decompensation

Aixia Guo^{1*}, Randi E. Foraker^{1,2}, Robert M. MacGregor³, Faraz M. Masood³, Brian P. Cupps³ and Michael K. Pasque³

¹ Institute for Informatics (^B), Washington University School of Modicine, St. Louis, MO, United States, ³ Department of Internal Medicine, Washington University School of Medicine, St. Lauis, MO, United States, ³ Department of Surgery, Washington University School of Medicine, St. Louis, MO, United States

Objective: Athough many clinical metrics are associated with proximity to decompensation in heart failure (HF), none are individually accurate enough to risk-stratity HF patients on a patient-by-patient basis. The dire consequences of this inaccuracy in risk stratification have profoundly lowered the clinical threshold for application of high-risk surgical intervention, such as ventricular assist device placement. Machine learning can detect non-intuitive classifier patterns that allow for innovative combination of patient feature predictive capability. A machine learning-based clinical tool to identify proximity to classtrophic HF deterioration on a patient-specific basis would enable more efficient direction of high-risk surgical intervention to those patients who have the most to gain from it, while sparing others. Synthetic electronic health information, and can be analyzed as if they were original data but without any privacy concerns. We demonstrate that synthetic EHR data can be easily accessed and analyzed and are amenable to machine learning analyses.

Methods: We developed synthetic data from EHR data of 26,575 HF patients admitted to a single institution during the decade ending on 12/31/2018. Twenty-seven clinically-relevant features were synthesized and utilized in supervised deep learning and machine learning algorithms (i.e., deep neural networks [DNN], random forest [RF], and logistic regression [LFI] to explore their ability to predict 1-year mortality by five-fold cross validation methods. We conducted analyses leveraging features from prior to/at and after/at the time of HF diagnosis.

Results: The area under the receiver operating curve (AUC) was used to evaluate the performance of the three models: the mean AUC was 0.80 for DNN, 0.72 for RF, and 0.74 for LR. Age, creatinine, body mass index, and blood pressure levels were especially important features in predicting death within 1-year among HF patients.



Decompensation

OPEN ACCESS

Huazhong University of Science and

Edited by:

Juan Liu.

Technology, China

Reviewed by:

Liang Zhang,

United States Konatao Chen.

Linited States

Aixia Guo aixia.guo@wustl.edu

Citation

*Correspondence:

Specialty section:

Health Informatics

a section of the journal

Frontiers in Digital Health Received: 27 June 2020

Accepted: 13 November 2020

Published: 07 December 2020

Pasque MK (2020) The Use of

Synthetic Electronic Health Record Data and Deep Leaming to Improve

Timing of High-Risk Heart Failure

Surgical Intervention by Predicting Proximity to Catastrophic

Front. Digit. Health 2:576945. doi: 10.3389/fdath.2020.576945

Guo A, Foraker RE, MacGregor RM, Masood FM, Cupps BP and

This article was submitted to

Xidian University, China Zhibo Wana.

University of Central Florida,

Liniversity of Pennsylvenia

December 2020 Volume 2 Article 576945



@TheInstituteDH #MEDINF023

RSEW SIDS

 We used a traditional adversarial approach to assess the privacy preserving nature of synthetic data as compared to de-identified data.

 Our results indicate that synthetic data cannot be confidently re-identified to the same level as deidentified data.

#MEDINF023



• We measured fingerprinting scores between sets of real and synthetic data to measure the amount of variability introduced via synthetic data.

• Measured the amount of information that could reliably be learned about real data from synthetic data.



@TheInstituteDH #MEDINF023

RSEW SIDS Aan



Postoperative Pregabalin Use







NEXT STEPS

ĠĊŗ ġĊŗ IMPORTANCE



ję.

IMPORTANCE



Key Takeaway: Need to be more thoughtful in implementing decision support, and it is helpful to measure patient impact

Next Step: Work with decision support team and clinicians



FINDINGS







Radiocontrast Switching in ED







IMPORTANCE

NEXT STEPS

କୁଚ୍ଚୁ କୁଚ୍ଚୁ କୁଚ୍ଚୁ



IMPORTANCE

NEXT STEPS

Key Takeaway: The better option for patients and staff is also better for the bottom line Next Step: Work with clinical leaders to implement, monitor benefits



FINDINGS



Causes of Hospitalizations in Ottawa







2003/4-2010/11 2011/12-2018/19

\$5M

\$0M

\$14K

2003/4-2010/11 2011/12-2018/19

\$16K

500

S

N

Ζ

උදුසු





Extend analysis to determine impact of EMR, covid-19, and other factors

Assess other diseases



Further analysis to determine possible explanations for decreasing cases and lung cancer trends

Focused efforts within FSAs to address possible SDOH \rightarrow impact of screening, smoking behaviour, other environmental risks



FINDINGS

222

STEPS

IMPORTANCE