Health Data Science Group

PUTTING MACHINE LEARNING INTO REAL-WORLD PRACTICE:

PATIENT-LEVEL PREDICTION DEVELOPMENT AND VALIDATION IN OBSERVATIONAL DATA



Peter R. Rijnbeek Associate Professor Health Data Science Department of Medical Informatics www.healthdatascience.nl



Problem definition





Among a target population (T), we aim to predict which patients at a defined moment in time (t=0) will experience some outcome (O) during a time-at-risk Prediction is done using only information about the patients in an observation window prior to that moment in time.



Model Development Pipeline



Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD) statement¹.

- Sharing of model development details
- Discrimination and Calibration
- Internal and external validation
- *Etc.*

Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD): Explanation and Elaboration

Karel G.M. Moons, PhD; Douglas G. Altman, DSc; Johannes B. Reitsma, MD, PhD; John P.A. Ioannidis, MD, DSc; Petra Macaskill, PhD; Ewout W. Steyerberg, PhD; Andrew J. Vickers, PhD; David F. Ransohoff, MD; and Gary S. Collins, PhD

The TBPCD (Transparent Reporting of a multivariable prediction model for Individual Prognosis of Distancem includes a 22-ben checklist, which aims to improve the reporting of studendopping, uddening, or updating a prediction model, prediction model tauly regardless of the study methods used. This explanation and elaboration document describes the ratioals; clarifishe meaning of ach there and discusses why transnaic clarifishe meaning of study there and discusses why transnaic clarifishe there and a study regardless of the study methods used. This explanation and elaboration document describes the ratioand clinical used under the production model. Each checklise them of the TBPCD Statement is explained in datafal and accompanied by published examples of good reporting. The document also provides a valuable reference of issues to consider when designing, conducting, and analyzing prediction model studies. To aid the editorial process and help peer reviewers and, ultimately, readers and systematic reviewers of prediction model studies, its recommended that authors include a completed checklist in their submission. The TBPOD checklist can also be downloaded from www.tripod.statement.org.

Ann Intern Med. 2015;162:W1-W73. doi:10.7326/M14-0698 www.annals.org For author affiliations, see end of text. For members of the TRPOD Group, see the Appendix.



¹ Moons, KG et al. Ann Intern Med. 2015;162(1):Wasked as internal/staff & contractors by the European Medicines Agency

Current status of predictive modelling:



Review of 422 papers with 579 models

External validation
 Internal validation
 No validation

Preliminary results of review led by Gynthia Yang of Erasmus MC



Current status of predictive modelling:



Preliminary results of review led by Oynthia Yang of Erasmus MC

Erasmus MC

What is needed?

Full transparency and reproducibility

- 1. Standardised Health Data with respect to structure and terminology
- 2. Standardised Analytical Pipelines that enforce best modelling practices
- 3. Share models and allow extensive external validation across many databases
- 4. Disseminate all performance results



Our mission for Patient-Level Prediction

The Observational Health Data Sciences and Informatics (OHDSI) and the European Health Data and Evidence Network (EHDEN) developed a systematic process to learn and evaluate large-scale patient-level prediction models using observational health data in a large data network



Patient-Level Prediction Framework



R-package

Design and implementation of a standardized framework to generate and evaluate patient-level prediction models using observational healthcare data \hat{o}

Jenna M Reps 🖾, Martijn J Schuemie, Marc A Suchard, Patrick B Ryan, Peter R Rijnbeek

Journal of the American Medical Informatics Association, Volume 25, Issue 8, August 2018, Pages 969–975, https://doi.org/10.1093/jamia/ocy032 Published: 27 April 2018 Article history ▼

🔎 PDF 🛛 💵 Split View 🛛 💪 Cite 🎤 Permissions 🛛 < Share 🔻

Abstract

Objective

To develop a conceptual prediction model framework containing standardized steps and describe the corresponding open-source software developed to consistently implement the framework across computational environments and observational healthcare databases to enable model sharing and reproducibility.

www.github.com/OHDSI/PatientLevelPrediction

Book-of-OHDSI

https://book.ohdsi.org/



EHDEN Academy https://academy.ehden.eu

Study Results https://data.ohdsi.org



Prediction team in action









Training

Extraction

Definition



Validation

Validation



Model Specification

ATLAS

Data Sources

VIEW:

CO Dros

and the second	
Dest	
2,46	
Alunia	
	D-L-

+ Add Model Se

Previous

Previous

Filter:

Filter:

(+4 more covariate settings)

Q Search								
☐ Concept Sets ☑ Chorts ☑ Chorts								
 Cohort Definitions Characterizations 	Show 10	▼ entries	M	Model Settings				
📥 Cohort Pathways	Remove	Name	- 17	Show 10	• entries			
7 Incidence Rates	×	[EHDEN RA] Female new users of methoxtrexate monotherapy used for PLP		51104 10				
Profiles	×	[EHDEN RA] New users of methoxtrexate monotherapy used for PLP	_ .	Remove	Model		Options	
Estimation	Showing	1 to 2 of 2 entries		×	LassoLogisticRegressionSettings		{"variance":0.01,"seed":null}	
Prediction	ant custom	- Chat	- :	Showing 1	to 1 of 1 entries			
Jobs	* Outco	me Cohorts						
Configuration	Show 10	▼ entries	C	ovariate S	Settings			
Feedback	Remove	Name		Column vi	isibility Copy CSV Show 10 T entries			
	×	[EHDEN RA] Leukopenia events using diagnoses and measurements	- 17	Remove	Options			
	×	[EHDEN RA] Pancytopenia events using diagnoses and measurements			DemographicsGender, DemographicsAge, Demograph	nicsAgeGroup DemographicsIndevMr	anth ConditionOccurrenceAnvTimePrior DrugGroupFral ongTerm (+4	
	×	[EHDEN RA] Pancytopenia or leukopenia events using diagnoses and measureme	nts		bemographicsdender, bemographicsAge, bemograph	incongeoroup, bemographicandexind		
	×	[EHDEN RA] Stroke (ischemic or hemorrhagic) events (any visit) (1)		×	DemographicsGender, DemographicsAgeGroup, Dem	ographicsIndexMonth, ConditionOccu	urrenceAnyTimePrior, CharlsonIndex, Dcsi (+1 more covariate settings)	
Apache 2.0 open source software	×	[EHDEN RA] Opportunistic Infections (2)		Showing 1	to 2 of 2 entries			
provided by	×	[EHDEN RA] Serious Infection events (2)						
OHDSI	×	[EHDEN RA] Serious Infection, opportunistic infections and other infections of int	ares Po	opulation	Settings			

Population	Settings					+ Add Population S
Column v	visibility Copy CSV	Show 10 🔻 entries				Filter:
Remove	Risk Window Start	Risk Window End	Washout Period	Include All Outcomes	Remove Subjects With Prior Outcome	🍦 Minimum Time At Risk
×	1d from cohort start date	730d from cohort start date	365d	true	true	1d
×	1d from cohort start date	90d from cohort start date	365d	true	true	1d
×	365d from cohort start date	1826d from cohort start date	365d	true	true	1d
						Dravioure
						T zajus

Generate R-Package and share with the world

😵 Review & Download								
Review Full Study Specification								
Please review the full study specification below and scroll down the page to download the study package.								
Full Analysis List 132 Prediction Problem Settings (2) Analysis Settings (6)								
	Target Cohort Name	Outcome Cohort Name	Model Name	Model Settings	Covariate Settings	Risk Window	Risk Window	
▼ Target Cohorts						Start	End	
(EHDEN RA) New users of methoxtrexate monotherapy used for PLP (66) [EHDEN RA) Female new users of methoxtrexate monotherapy	[EHDEN RA] New users of methoxtrexate monotherapy used for PLP	[EHDEN RA] Serious Infection, opportunistic infections and other infections of interest event (1)	LassoLogisticRegressionSettings	{"variance":0.01,"seed":null}	"attr_class":"covariateSettin	1	730	
Outcome Cohorts [EHDEN RA] Leukopenia events using diagnoses and measurements (12) [EHDEN RA] Consecturistic	[EHDEN RA] New users of methoxtrexate monotherapy used for PLP	[EHDEN RA] Serious Infection, opportunistic infections and other infections of interest event (1)	LassoLogisticRegressionSettings	{"variance":0.01,"seed":null}	"attr_class":"covariateSettin	1	90	
T Model Name LassoLogisticRegressionSettings (132)	[EHDEN RA] New users of methoxtrexate monotherapy used for PLP	[EHDEN RA] Serious Infection, opportunistic infections and other infections of interest event (1)	LassoLogisticRegressionSettings	{"variance":0.01,"seed":null}	"attr_class":"covariateSettin	365	1826	
▼ Risk Window	[EHDEN RA] New users	[EHDEN RA] Serious Infection,						





Prediction Viewer About Internal Validation External Validation Plot Table Evaluation Summary Characterization ROC Calibration Demographics Preference Box Plot Settings Train Test Included ROC Plot ROC Plot Not included outcome Line 0.8 0.8 0.8 Prevalance in persons with 0.6 Sensitivity Sensitivity 0.4 0 0.4 0.2 0. 0.2 0.2 0.4 0.6 0.8 0.2 0.4 0.6 0.8 1 0.2 0.4 0.6 0.8 0 Prevalance in persons without outcome 1-specificity 1-specificity

Figure 2. The Receiver Operating Characteristics (ROC) curve shows the ability of the model to discriminate between people with and without the outcome during the time at risk. It is a plot of sensitivity vs 1-specificity at every probability threshold. The higher the area under the ROC plot the higher the discriminative performance of the model. The diagonal refers to a model assigning a class at random (area under de ROC = 0.5).

Figure 1. The variable scatter plot shows the mean covariate value for the people with the outcome against the mean covariate value for the people without the outcome. The meaning of the size and color of the dots depends on the settings on the left of the figure.

Data generated: 2018-10-01 01:11:16

Share model performance



Erasmus MC

zalm

Large scale validation and dissemination

The tool auto generates a word document containing all the model specifications, internal and external validation results, model details etc. etc. which serves as a kickstart for result dissemination.

Multiple interesting visualisation can be created:



Seek COVER: COVID risk prediction

Objective: develop and externally validate **COV**ID-19 **E**stimated **R**isk scores that quantify a patient's risk of hospital admission, hospitalization requiring intensive services or fatality.



doi: https://doi.org/10.1101/2020.05.26.20112649

- 14 data sources from 6 countries
- Externally validated in 44,507 COVID cases from 5 data sources in South Korea, Spain, USA



Classified as internal/staff & contractors by ton depreview dicines Agency

Results dissemination





What is next: Prediction Model Library





zalus

We need external validation at scale!

JMIR Med Inform. 2021 Apr; 9(4): e21547. Published online 2021 Apr 5. doi: 10.2196/21547 PMCID: PMC8023380 PMID: 33661754

Implementation of the COVID-19 Vulnerability Index Across an International Network of Health Care Data Sets: Collaborative External Validation Study

Monitoring Editor: Christian Lovis

Reviewed by David Maslove, JianLi Wang, and Anoop Austin

Janna M. Rens, BSc, MSc, PhD,⁸¹ Chungsoo Kim, PharmD,² Ross D. Williams, MSc,³ Aniek F. Markus, BSc, MSc,³ Cynthia Yano, BSc, MSc,³ Taita Duarte-Salles, MPH, PhD,⁴ Thomas Falconer, BSc, MSc,⁵ Jittendra Jonnagadala, MIS, PhD,⁶ Andrew Williams, PhD,⁷ Sergio Famidace-Bertolin, MSc,⁴ Sout, LU/Vall, PhD,⁸ (ristin Koasta, MPH,⁹ Gowtham Rao, MD, PhD,¹ Azza Shoaibi, PhD,¹ Anna Ostrozolets, MD,⁵ Matthew E Spotnitz, MPH, MD,⁵ Lin Zhang, PhD,^{10,11} Paula Casajiust, BSc, MSc,¹2 Ewoolt W Steverberg, MSc, PhD,^{13,14} Fredrik Nyberg, MPH, MD, PhD,¹⁵ Benjamin Skov Kasa-Hansen, MSc, MD,^{16,17} Young Hwa Chol, MD, PhD,¹⁹ Daniel Mornles, PhD, MBChB,¹⁹ Siaw-Teng Liaw, PhD,⁶ Maria Tereza Fermandes Abraha, PhD,²⁰ Carlos Areis, MSc, PT,²¹ Michael E Matheny, MD, MPH, MS,²² Kristine E Lynch, PhD,⁶ Maria Araoón, MSc,⁴ Rae Woong Park, MD, PhD,²² George Hripcsak, MD, MS,⁵ Christian G Reich, MD, PhD,⁹ Maria Asubard, MD, PhD,²⁴ Seng Chan You, MD, MS,²³ Patrick B Ryan, PhD,¹ Daniel Priot-Ahambra, MD, PhD,⁵ and Peter R Ripbeek, PhD³

Join the network!

Pope at al. BMC Madical De

Reps et al. BMC Medical Research Methodology (2020) 20:102 https://doi.org/10.1186/s12874-020-00991-3 BMC Medical Research Methodology

Open Access

RESEARCH ARTICLE

Feasibility and evaluation of a large-scale external validation approach for patientlevel prediction in an international data network: validation of models predicting stroke in female patients newly diagnosed with atrial fibrillation

Jenna M. Reps^{1*}, Ross D. Williams², Seng Chan You³, Thomas Falconer⁴, Evan Minty⁵, Alison Callahan⁶, Patrick B. Ryan¹, Rae Woong Park^{3,2}, Hong-Seok Lim⁶ and Peter Rijnbeek²

Abstract

Background: To demonstrate how the Observational Healthcare Data Science and Informatics (DHDS) collaborative network and standardization can be utilized to scale-up external validation of patient-level prediction models by enabling validation across a large number of heterogeneous observational healthcare datasets.

Methods: Five previously published prognostic models (ATRA CHADS), CHADS/MASC, D-Stoeke and Framingham that predict future tisk of stroke in patients with anial fibriliation were replicated using the CHOS frameworks. A network study was run that enabled the five models to be externally validated across nine observational healthcare datasets spanning three countries and five independent sizes.

Results: The five existing models were able to be integrated into the CHCI3 framework for patient-level prediction and they obtained mean catatitics canging between CS2-O63 across the 6 databases with sufficient data to predict stroke within 1 year of initial and Britliation diagnosis for females with atrial fibrillation. This was comparable with existing validation studies. The validation network study was published at https://github.com/OHDS/ Sub/phrotoc32mbot/refmarkef2mbitingforvbef3bc4acreanValidation.

(Continued on next page)

Using the OHDSI network to develop and externally validate a patient-level prediction model for Heart Failure in Type II Diabetes Meilitus

📀 Ross D. Williams, Jenna M. Reps, Jan A Kors, Patrick B Ryan, Ewout Steyerberg, Katia M. Verhamme, Peter R. Rijnbeek

doi: https://doi.org/10.1101/2021.04.06.21254966

This article is a proprint and has not been peer-reviewed [what does this mean?]. It reports new medical research that has yet to be evaluated and so should not be used to guide clinical practice.

Abstract Full Text Info/History Metrics

Abstract

Introduction Heart Failure (HF) and Type 2 Diabetes Mellitus (T2DM) frequently coexist and exacerbate symptoms of each other. Treatments are available for T2DM that also provide beneficial treatment effects for HF. Guidelines recommend that patients with HF should be given Sodium-glucose oc-transporter-2 inhibitors in preference to other second-line treatments for T2DM. Increasing personalization of treatment means that patients who have or are at risk of HF receive a customised treatment. We aimed to develop and externally validate prediction models to predict the 1-year risk of incident HF in T2DM patients starting second-line treatment.



O Comment on this paper

Preview PDF

Acknowledgement

This work would not have been possible without many contributors in the OHDSI community and the EHDEN project







tion and the second funding from the language

innovative medicines initiative

This project has received funding from the Innovative Medicines Initiative 2 Joint Undertaking (JU) under grant agreement No 806968. The JU receives support from the European Union's Horizon 2020 research and innovation programme and EFPIA.

efpia

