

## Disclaimer

The views expressed here are the personal views of the authors and may not be understood or quoted as being made on behalf of or reflecting the position of the regulatory agency/agencies or organizations with which the authors are employed/affiliated.

- Ability to discriminate between different risk groups
- Improves patient outcomes by informing treatment decisions

## What is a good prediction model?

### General requirements

- Generates accurate predictions in individuals from potential population(s) for clinical use

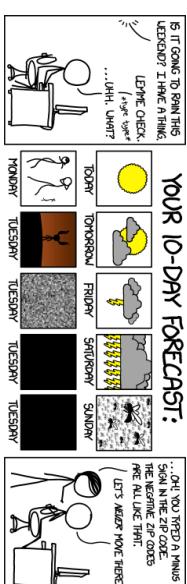


## Statistical methods for IPD-MA of prognosis studies



Valentijn de Jong<sup>1,2</sup>, PhD  
1. Julius Center for Health Sciences and Primary Care  
2. Data Analytics and Methods Task Force, European Medicines Agency

## What is a good prediction model?



<http://xkcd.com/1245/>



## The reality

**Most prediction models are not as good as we think**

- Quality of many prognostic model studies is poor
  - Absence of a study protocol
  - Exclusion of eligible study participants
  - Poor handling of missing data
  - Complex modelling in small samples
  - Incomplete registrations & reporting
- Internal validation too optimistic
- Lack of external validation



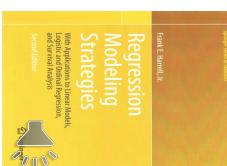
## Lack of external validity

Poor performance in new patients from different (but related) populations

### The reality

Poor performance in new patients from the development population

- Overfitting
- Inclusion of noise variables (e.g. via stepwise selection)
- Poor handling of missing values



## Numerous models for same target population + outcomes

RESEARCH

OPEN ACCESS

Prediction models for cardiovascular disease risk in the general population: systematic review

RESEARCH

FAST TRACK

Systematic review and critical appraisal

RESEARCH

OPEN ACCESS

Prediction models for diagnosis and prognosis of covid-19:

RESEARCH



Open Access

Check for updates

Editorial

Correspondence

Review

Supplementary material

Editorial

Commentary

Case report

Letter

Review

Case series

Editorial

Commentary

Case report

## Prediction models for COVID-19

731 models from 412 studies included in the review:

- **16 models** to identify subjects at risk in the general population
  - 75 based on medical images (deep learning)
- **593 prognosis models** for predicting mortality risk, progression to severe disease, or length of stay, or other
  - 75 based on medical images (deep learning)

## Prediction models for COVID-19

- Participants domain: **107/606 (18%) at high risk of bias**
  - Non-representative of the target population (e.g., non-consecutive patients)

- Predictors domain: **27/606 (4%) at high risk of bias**
  - Predictors not available at time of intended model use

- Outcome domain: **106/606 (17%) at high risk of bias**
  - Small sample size (->overfitting & no adjustment), incomplete reporting of model performance (e.g., no calibration)



## Prediction models for COVID-19

### Prediction models for diagnosis and prognosis of covid-19: systematic review and critical appraisal

#### FAST TRACK

Laure Wijnants,<sup>1,2</sup> Ben Van Calster,<sup>2,3</sup> Gary S. Collins,<sup>4,5</sup> Richard D'Heij,<sup>6</sup> Georg Henze,<sup>7</sup> Ewoud Schuit,<sup>1,2</sup> Elena Abu<sup>8</sup>, Ben Van Calster,<sup>2,3</sup> Amanesh Asfi,<sup>9</sup> Vanesa Belotti,<sup>10</sup> Marc M. Bonten,<sup>11</sup> Darren L. Daily,<sup>12</sup> Johanna A. Dammer,<sup>13</sup> Thomas P. De Bruyn,<sup>14</sup> Valentin M. de Jong,<sup>15</sup> Maarten De Vos,<sup>16</sup> Paula Dhama,<sup>17</sup> Sophie Eison,<sup>18</sup> Shan Gao,<sup>19</sup> Natacha Challe,<sup>20</sup> Michael O'Hartney,<sup>21</sup> Lieven Herckx,<sup>22</sup> Pauline Huygen,<sup>23</sup> Jeroen Hoogendoorn,<sup>24</sup> Mohamed Hudda,<sup>25</sup> Kevin Hemmers,<sup>26</sup> Michael Kammer,<sup>22</sup> Anna Lohmann,<sup>27</sup> Kim Luijten,<sup>28</sup> Ille Ma,<sup>29</sup> James B. Peeling,<sup>30</sup> Jack Wilkinson,<sup>31</sup> Arndur Nakano,<sup>32</sup> Karel G. M. Moons,<sup>33</sup> Maarten Van Smeden<sup>34</sup>

For unrecorded affiliations see end of the article

Correspondence: Dr. Laure Wijnants

Published online: 17 July 2020  
Received: 17 January 2020  
Accepted: 17 July 2020  
Editorial office: Dr. Laure Wijnants  
<http://doi.org/10.31233/osf.io/11x6cm>

© 2020 Wijnants et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction is free, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No permission is required from the copyright holders for this.

Final version accepted:

17 July 2020

Data sources

### New concept: Living Review

#### DATA EXTRACTION

Data supplement

Files in this Data Supplement:

Data Supplement - Web Appendix: Supplementary material

Original article DOI: Published 7 April 2020

Update article DOI: Published 21 July 2020

Update 2 article DOI: Published 21 July 2020

Population	Objective	AUC range
General population	Diagnosis or prognosis	0.52 to 0.99
Patients with covid-19	Prognosis	0.49 to 1.00



## The rise of big datasets

### The QR ESEARCH database

- Anonymised health records of over 25 million people from 1500 general practices spread throughout the UK
- Linkage to Hospital Episode Statistics, Mortality and Cancer Registration data



## Prediction models for COVID-19

- Many COVID-19 prediction models were poorly reported, at high risk of bias, and their reported performance is probably optimistic

- Application of any prediction models that was developed early on in the pandemic was not recommended

- Only 7 out of 606 had a low risk of bias.

## The rise of big datasets

Data increasingly available for thousands or even millions of patients from multiple practices, hospitals, or countries.

- Meta-analysis of individual participant data (IPD) from multiple studies
  - Observational studies
  - Randomized controlled trials
- Analyses of databases and registry data containing e-health records



## The rise of big datasets

### Why do we need big datasets?

- Development of better prediction models
- More extensive testing of model performance

OPEN ACCESS  
GUIDELINES AND GUIDANCE

### Individual Participant Data (IPD) Meta-analyses of Diagnostic and Prognostic Modeling Studies: Guidance on Their Use

Thomas P. A. Debray , Richard D. Riley, Maroeska M. Rovers, Johannes B. Reitsma, Karel G. M. Moons,  
Cochrane IPD Meta-analysis Methods Group 

Published: October 13, 2015 • <https://doi.org/10.1371/journal.pmed.1001886>

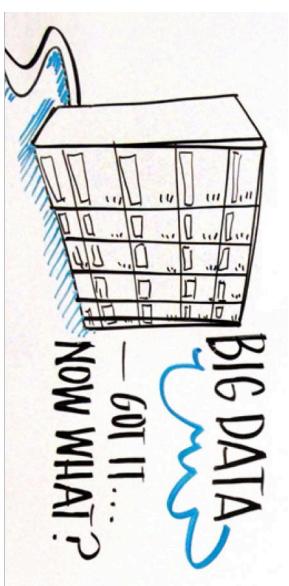


## External validation of prediction models in big datasets

### The rise of big datasets

### CALIBER

- EHR data encompassing more than 10 million adults with 400 million person-years of follow-up
- Primary care consultations and hospitalisations
- Clinical examination findings, blood laboratory results, prescriptions and vaccinations
- Diagnoses of diseases and mortality data



## Why do we need external validation?

- The predictive performance of a model estimated on the development data is often too optimistic
- A prognostic model should provide predictions that are valid outside the specific context of the sample that was used for model development
- How a model was derived is of little importance if it performs well.

## Measures of prediction model performance

### Overall performance

- Amount of explained variation ( $R^2$ )



Ref: Steyerberg. Clinical prediction models: a practical approach to development, validation and updating. Springer 2009.



## Background

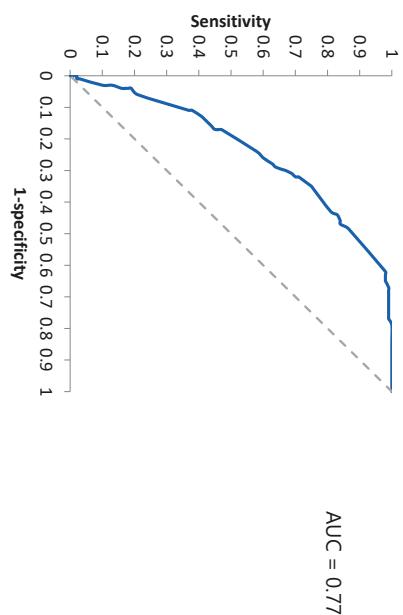
### Key references

- Riley et al. External validation of clinical prediction models using big datasets from e-health records or IPD meta-analysis: opportunities and challenges. *BMJ* 2016.
- Debray et al. Individual Participant Data (IPD) meta-analyses of diagnostic and prognostic modeling studies: guidance on their use. *PLOS MED* 2015.
- Evaluate the predictive accuracy
  - Overall performance
  - Calibration
  - Discrimination



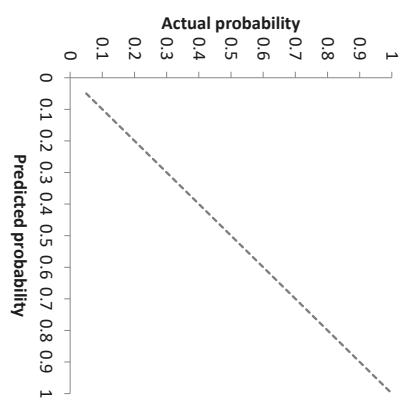
## Measures of prediction model performance

### Discrimination



## Measures of prediction model performance

### Calibration plot – good model?



## Measures of prediction model performance

### Discrimination

Quantifies the model's extent to distinguish between events and non-events

- Summary statistics
  - Concordance ( $C$ ) index
  - Area under the ROC curve (AUC)
  - Discrimination slope
- Visual inspection
  - Receiving Operating Characteristics (ROC) curve

### Calibration

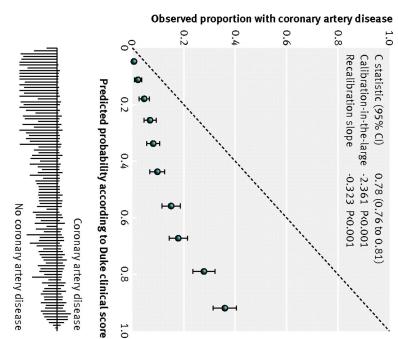
Agreement between observed outcomes and predictions

- Visual inspection
  - Calibration plot
- Summary statistics
  - O:E statistic (#observed events / #predicted events)
  - Calibration-in-the-large (or calibration intercept, which is related)
    - Calibration slope
  - Hosmer-Lemeshow goodness-of-fit test



## Measures of prediction model performance

Calibration plot – good model?

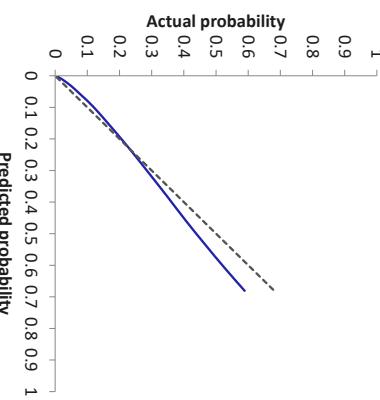


Ref: Genders et al. Prediction model to estimate presence of coronary artery disease: retrospective pooled analysis of existing cohorts. *BMJ*.



## Measures of prediction model performance

Calibration plot – good model?

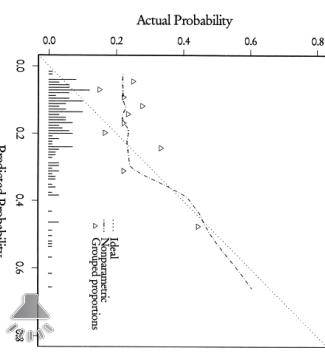


## Example study

### External validation of the model

- Population: 179 children from a different time period and another hospital from a different city in NL
- Very similar in- and exclusion criteria and variable def.
- AUC = 0.57 (95% CI: 0.47 – 0.67)

Ref: Bleeker SE et al. External validation is necessary in prediction research: a clinical example. *Journal of Clinical Epidemiology* 2003.



## Example study

Development and validation of a prediction model for the presence of serious bacterial infections in children with fever.

### Development of the model

- Population: 379 children between 1 month and 36 months of age referred to the Emergency Department from a hospital in NL (75 events)
- Analysis: logistic regression with forward stepwise variable selection (57 variables -> 9 predictors)
- Internal validation:
  - AUC = 0.825 (95% CI: 0.78 – 0.87)
  - Bootstrap-corrected AUC = 0.756



## Current shortcomings of validation studies

### Why do we need big datasets for external validation?

- The predictive performance of a model tends to vary across settings, populations and periods
- Multiple external validation studies are needed to fully appreciate the generalizability of a prediction model
- Heterogeneity in model performance is rarely assessed, but investigating its extent is crucial to evaluate the model's potential generalizability and clinical usefulness.



## Current shortcomings of validation studies

### External validation requires sufficient data

- Recommendations: > 100 events and > 100 non-events
- Less data available for model development
- Not all validation studies are equally informative.
  - To what extent do individuals from the validation sample represent the target population?
  - To what extent are estimates of model performance affected by flaws in the design and analysis of the validation study?
  - To what extent can the CPM be implemented across different populations and settings?



## Causes of heterogeneity in model performance

### Discrepancies in outcome and predictor assessment

- Different measurement method for predictors (e.g. using equipment from different manufacturers)
- Different recording time of predictors (e.g. before or after surgery)
- Different quantification of predictors (e.g. use of cut-points may vary)
- Different disease and outcome definitions
- Different follow-up lengths



## Causes of heterogeneity in model performance

### Invalid predictor effects

- Over-fitting of the prediction model to the development study (sometimes avoided using *penalization*)
  - Biased estimates of predictor effects (e.g. due to flaws in the development study)
  - Missed interactions or non-linear associations



## Causes of heterogeneity in model performance

### Case-mix variation (spectrum effect)

- Different distribution of predictor values
  - Different standards of care and treatment strategies
  - Different starting points  
(e.g. earlier diagnosis due to screening program)
  - Different outcome prevalence or incidence
  - Different participant or setting characteristics

Case-mix variation can lead to genuine differences in the performance of a prediction model, even when the predictor effects remain "correct" in the validation study



## Causes of heterogeneity in model performance

## Differences between study characteristics

# BMJ Open Empirical evidence of the impact of study characteristics on the performance of prediction models: a meta-epidemiological study

Johanna A A G Damen,<sup>1,2</sup> Thomas P A Debray,<sup>1,2</sup> Romin Pajouheshnia,<sup>2</sup> Johannes B Reitsma,<sup>1,2</sup> Rob J P M Scholten,<sup>1,2</sup> Karel G M Moons,<sup>1,2</sup> Lotty Hooft<sup>1,2</sup>



## Interpretation of model performance

Need to adjust for clustering

Does ignoring clustering in multicenter data influence the performance of prediction models? A simulation study

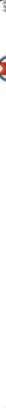
Presently, no comprehensive framework for the evaluation of model performance on multicenter datasets is available. In this article, we propose a framework for the evaluation of model performance on multicenter datasets. At the center of our approach is the concept of model validation, which is based on the choice of the model (standardized main effects, logistic regression), the level of prediction (marginal versus conditional) and the level of model validation (center versus population). In particular, when data is heavily clustered (ICC > 0.5), center-specific predictions offer the best predictive performance at the model level and the center level. We recommend that models should reflect the data structure, while the level of model validation should reflect the research question.

**SMMR**  
STATISTICAL METHODS IN MEDICAL RESEARCH

## Interpretation of model performance

Need to disentangle case-mix differences from differences in predictor effects!

 ELSEVIER



*Journal of Clinical Epidemiology* 66 (2013) 279–289

**ORIGINAL ARTICLES**

A new framework to enhance the interpretation of external validation studies of clinical prediction models

Journal of Clinical Epidemiology 68 (2015) 279–289  
www.jclinepi.com  
A new framework to enhance the interpretation of external validity studies of clinical prediction models  
Thomas P.A. Debray<sup>a,b</sup>, Yvonne Vergauwe<sup>b</sup>, Hendrik Koffijberg<sup>b</sup>, Daan Nijs<sup>b</sup>,  
Evert W. Steyerberg<sup>a</sup>, Karel G.M. Moons<sup>a,c</sup>  
<sup>a</sup>Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, PO Box 8550,  
Utrecht, The Netherlands  
<sup>b</sup>Department of Public Health, Erasmus University Rotterdam, The Netherlands  
Accepted 30 June 2014; published online 18 August 2014



## Examining heterogeneity and improving model performance

### Guidance paper

RESEARCH METHODS AND REPORTING



A guide to systematic review and meta-analysis of prediction model performance  
Thomas P. Debray,<sup>1,2</sup> Johanna A.G. Dammen,<sup>1,2</sup> Koen I.F. Stell,<sup>3</sup> Joie Ensor,<sup>3</sup> Lotte Hooft,<sup>1,2</sup> Johannes B. Reitsma,<sup>1,2</sup> Richard D. Riley,<sup>2</sup> Karel G.M. Moons,<sup>1,2</sup>

- Previously developed prediction model for diagnosing DVT
- Logistic regression analysis
- Three predictors
  - Sex
  - Surgery
  - Calf difference

$$Pr(DVT) = \frac{1}{1 + \exp(-(\alpha + \beta_1 \text{sex} + \beta_2 \text{surg} + \beta_3 \text{cdif}))}$$



## Examining heterogeneity and improving model performance

### Recommendations

- Calculate key performance statistics in each cluster (e.g. study or hospital)
- Summarize the performance measures by applying (multivariate) random effects meta-analysis
- Quantify between-study heterogeneity in model performance using 95% prediction intervals

### More guidance

- Prognostic model research
- Systematic reviews and meta-analysis of prognosis research studies
  - Individual participant data meta-analysis of prognosis studies
  - Electronic health care records and prognosis research

OXFORD

PROGNOSIS RESEARCH  
IN HEALTHCARE

Concepts, Methods, and Impact

EDITED BY  
Richard D. Riley • Danielle A. van der Windt  
Peter Croft • Karel G.M. Moons

## Example 1 (same as practical)

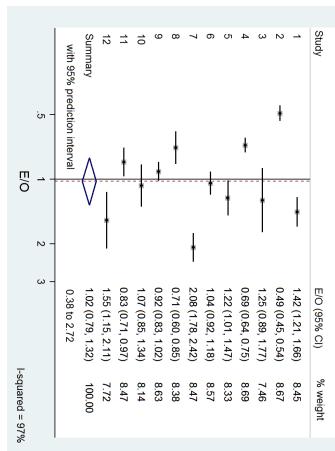
### Diagnosis of deep vein thrombosis in patients suspected of DVT

Geersing GJ, Zutphoff NPA, Kearon C, Anderson DR, Cate-Hoek AJ ten, Eijf JL, et al. Exclusion of deep vein thrombosis using the Wells rule in clinically important subgroups: individual patient data meta-analysis. *BMJ*. 2014 Mar 10;348:b1340. Available from: <https://www.bmjjournals.org/content/348/bmj.g1340>



## Example 1

### External validation in 12 studies



I-squared = 97%

- 364 practices from the UK
- Total sample size N=2,084,445
- Total number of events E=93,564

Again, each cluster might be viewed as a different external validation study!

## Example 2

### Prognosis of cardiovascular disease in patients from general practices using QRISK<2

Hippisley-Cox J, Coupland C, Vinogradova Y, Robson J, Minhas R, Sheikh A, et al. Predicting cardiovascular risk in England and Wales: prospective derivation and validation of QRISK2. *BMJ*. 2008 Jun 28;336(7559):1475-82. Available from: <https://doi.org/10.1136/bmjj.360944967625>



I-squared = 97%

## Example 1

### External validation in 12 studies

- Sample size: 153 – 1768 (total N=10014)
- Event occurrence: 8% - 39% (total E=1897)

- Results (95% confidence interval)
  - Calibration-in-the-large: -0.004 (-0.313; 0.305)
  - Calibration slope: 0.980 (0.853; 1.107)
  - E/O ratio: 1.02 (0.81; 1.28)
  - C-statistic: 0.687 (0.669; 0.705)

Does the model predict well? What about generalizability?

## Example 1

### External validation in 12 studies

- Substantial between-study heterogeneity!
- Approximate 95% prediction intervals:
  - Calibration slope = 0.59 to 1.38
  - E/O ratio = 0.38 to 2.72
  - c-statistic = 0.64 to 0.73

The model requires improvements to improve discrimination and to be clinically useful!

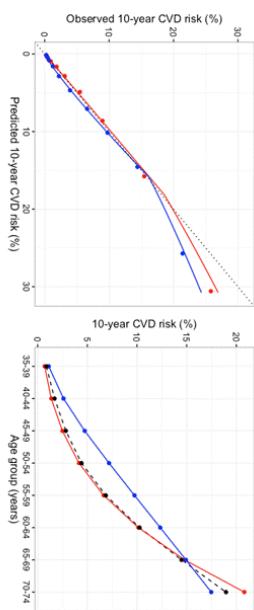
## Example 2

### External validation in 364 practices

## Example 2

### Investigating between-study heterogeneity

- Recall that variation in case-mix severity and case-mix heterogeneity may affect model performance.
  - Larger case-mix variation is related to larger discrimination performance
  - Populations with a narrower case-mix tend to have worse discrimination performance

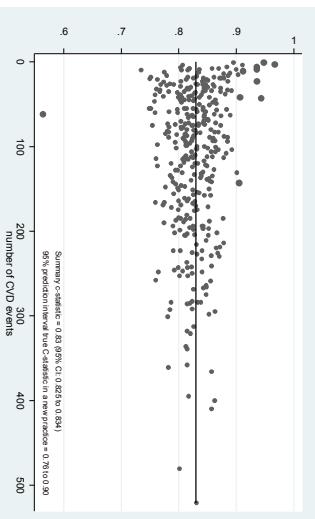


## Example 2

### External validation in 364 practices

## Example 2

### External validation in 364 practices



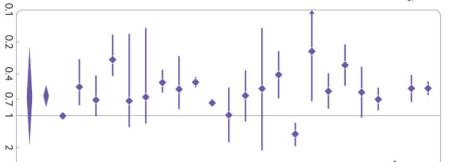
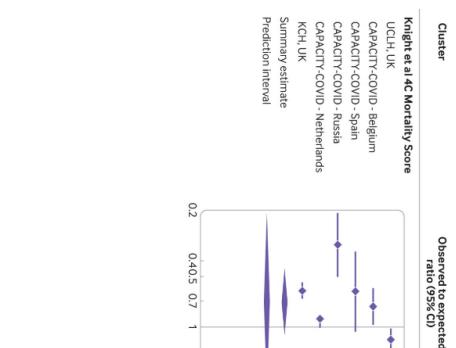
### External validation in 364 practices

- Discrimination
  - Summary c-statistic: 0.83
  - 95% confidence interval: 0.825 to 0.834
  - 95% prediction interval: 0.76 to 0.90
- Calibration
  - E/O summary estimate: 1.01
  - Slight over-prediction in women at higher CVD risk
  - QRISK2 appears to accurately predict 10-year CVD risk across all age groups.

## Example 3: prognosis in COVID-19 patients

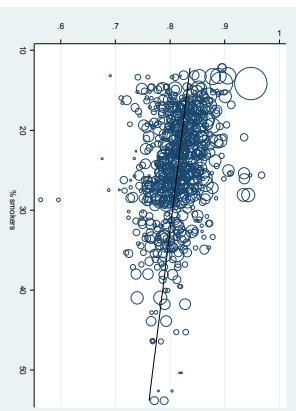
- External validation in 46,914 patients from 18 countries
- Design: Two stage individual participant data meta-analysis.
- Setting: Secondary and tertiary care.
- Patients: diagnosed with covid-19
- Outcome: 30-day mortality or in-hospital mortality.

De Jong VMT, Rousset RZ, Antonic-Villa NF, Bueno AG, Calster BV, Bello-Chavolla OY, et al. Clinical prediction models for mortality in patients with covid-19 - external validation and individual participant data meta-analysis. *BMJ* 2022 Jul 12;370:e69801. Available from: <https://doi.org/10.1136/bmj-2021-063881>



## Example 2

### Investigating between-study heterogeneity



NB Circle size is weighted by the precision of the c-statistic estimate (i.e. larger circles indicate c-statistic estimates with smaller standard errors, and thus more weight in the meta-regression)

## Example 3: prognosis in COVID-19 patients

### Cluster

#### Knight et al 4C Mortality Score

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

Wang et al clinical model

C statistic  
(95% CI)

Cluster

## Development of prediction models in big datasets



## Model development using big datasets



### Practical and methodological challenges

Caution is warranted when interpreting summary estimates of model performance and between-study heterogeneity.

- **Data quality**
  - Missing predictor values
  - Non-standardised definitions of diagnoses and outcomes
  - Incomplete follow-up times and event dates
  - Lack of recording of novel/costly predictors
  - Risk of double entries
- **Data dredging**

Need for study protocols and quality appraisal tools!



Need for study protocols and quality appraisal tools!



## Model development using big datasets

So, let's pool our IPD and 'launch' the analysis?



## Current practice

Abstract of BMC Medical Research Methodology 2014, 14:4  
<http://www.biomedcentral.com/1471-2288/14/4>



RESEARCH ARTICLE

Open Access

### Developing and validating risk prediction models in an individual participant data meta-analysis

Rahel Ahmed<sup>1</sup>, Thomas P. Delaney<sup>2</sup>, Karl G. Moons<sup>3</sup> and Richard D. Riley<sup>4\*</sup>

**Abstract** Risk prediction models estimate the risk of developing future outcomes for individuals based on one or more underlying characteristics (predictors). We review how researchers develop and validate risk prediction models within an individual participant data (IPD) meta-analysis, and reporting in 15 articles that developed a risk prediction model using IPD from multiple sources.

**Method** A qualitative review of the 15 articles that developed a risk prediction model using IPD from multiple sources.

**Results** The IPD approach offers many opportunities but methodological challenges exist, including unavailability of outcome data, missing participant data and incomplete data. Most studies used a cross-validation approach to validate their models, but did not allow for any study differences in baseline risk (stratification), potentially limiting their models' applicability and performance. In some publications, only two articles used external validation on different data, including a novel method which compares their results with those of the IPD source, test performance in the outcome study and repeats by others.

**Conclusions** An IPD meta-analysis offers unique opportunities for risk prediction research. Researchers can make more of this by allowing separate model intercept terms for each study population to improve generalisability, and by using a hierarchical cross-validation to simultaneously develop and validate their model. Methodological challenges can be reduced by prospectively planned collaborations that share IPD and risk prediction.

**Keywords:** Meta-analysis, prognostic factor, prognosis, individual participant (patient) data, flowchart, reporting

## Model development using big datasets

### Main opportunities

- Increase total sample size
- Increase available case-mix variability
- Ability to standardize analysis methods across IPD sets
- Ability to directly validate developed prediction models across a wide range of populations and settings
- Ability to evaluate generalizability of the model

## Model development using big datasets

### Wait... which analysis?



## Current practice

### Investigation of heterogeneity

- 12 articles did not consider heterogeneity in predictor effects
  - 1 article investigated interaction terms between study and each predictor
  - 1 article investigated heterogeneity using the  $\chi^2$  statistic
  - 1 article investigated heterogeneity using Chi-square test



## Problems with meta-analysis methods

### Random effects summaries are of limited value

- Predictor effects and/or baseline risk may take different values for each included study
  - Which parameters to use when validating/implementing the model in new individuals or study populations?
  - When do study populations differ too much to combine?

Need for a framework that can identify the extent to which aggregation of IPD is justifiable, and provide the optimal approach to achieve this.



## Current practice

### Analysis methods (review of 15 IPD-MA)

- 10 articles pooled all the IPD into one big dataset and analysed it ignoring clustering of patients
  - 4 articles used a one-stage approach accounting for clustering (e.g. Stratification of intercept term)
  - 1 article used a two-stage approach accounting for clustering



## Recommendations

- **Allow for different baseline risks in each of the IPD studies**
  - Account for differences in outcome prevalence (or incidence) across studies
  - Examine between-study heterogeneity in predictor effects and prioritize inclusion of (weakly) homogeneous predictors
- **Implement a framework that uses internal-external cross-validation**

## The framework

**Step 2:** Choosing an appropriate model intercept when implementing the model to new individuals

- Average intercept term (e.g., pooled estimate)
- Updating of intercept term (requires patient-level data)
- Use intercept of included study (e.g., based on outcome occurrence)

Propose which intercept term to use in new populations  
!! More difficult in case of heterogeneous predictor effects



## Problems with meta-analysis methods

### The framework

#### Step 1: Different choices for combining IPD

- Merge all data into one big dataset and ignore heterogeneity
- Allow heterogeneous baseline risk across studies\*
  - assume random effects distribution for the intercept terms
  - estimate study-specific intercept terms
- Advanced modeling of predictor effects is also possible
  - Nonlinear effects
  - Interaction terms



Research Article

### A framework for developing, implementing, and evaluating clinical prediction models in an individual participant data meta-analysis

Thomas P.A. Deetray Karel G.M. Moons, Ikhlaaq Ahmed, Hendrik Koltijberg,

Richard David Riley

First published: 11 January 2013 Full publication history

DOI: 10.1002/sim.5732 View/Save citation

Cited by: 19 articles Refresh Cite this article



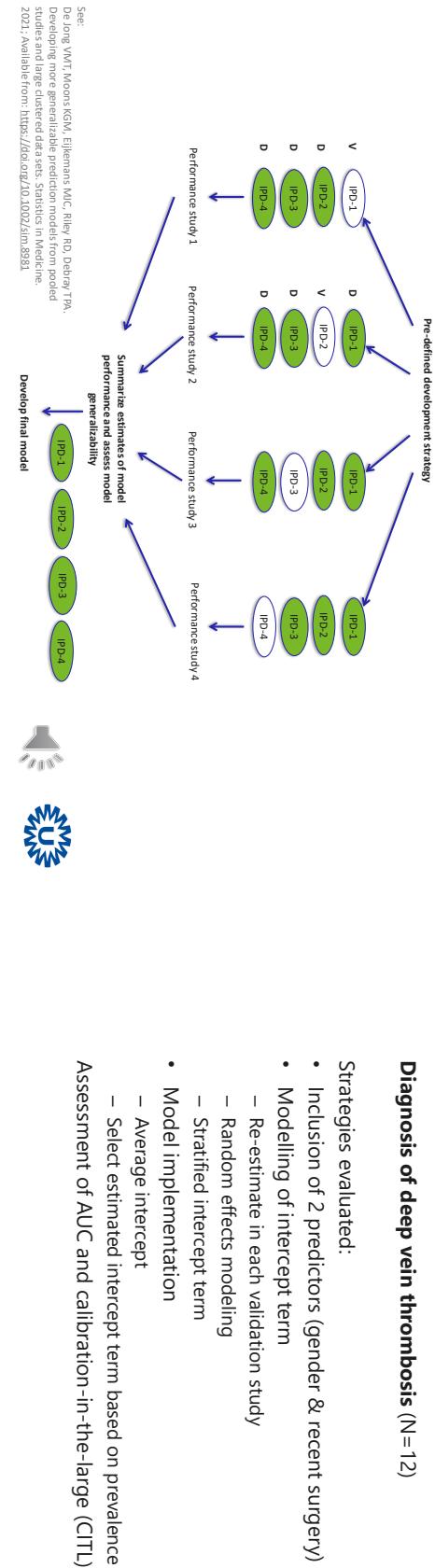
\*For info on this, see Faloutsos N, Galter BV, Timmerman D, Wyant L. Developing risk models for multicenter data using semiparametric regression produced suboptimal predictions: A simulation study. *Biometrical Journal*. 2020;62(4):932–44. Available from: <https://onlinelibrary.wiley.com/doi/10.1002/bimj.201900075>



## Internal-external cross-validation

### Example 1

#### Diagnosis of deep vein thrombosis (N=12)



## The framework

### Step 3: Model evaluation to check whether...

- Strategy for estimating predictors and intercept is adequate
  - Strategy for choosing intercept term (and predictor effects) in new study population is adequate
  - Model performance is consistently well across studies
    - Discrimination
    - Calibration
- => Use of internal-external cross-validation

## Formally comparing different strategies

### • Meta-analyze estimates of model performance

Compare summary estimates

– Compare prediction intervals

### • Rank different strategies by their overall performance

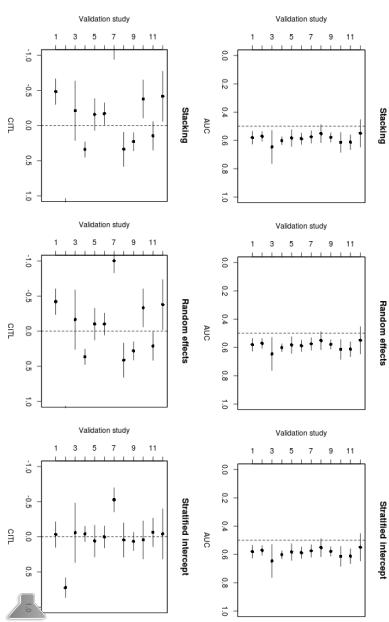
– Calculate the joint probability that, in a new population, model performance will meet certain criteria  
(e.g. C-statistic > 0.7 and cal. slope between 0.9 and 1.1)

– The strategy with the largest probability will be ranked first  
– This requires reliable prediction intervals!

## Example 1

### Example 1

#### Diagnosis of deep vein thrombosis (N=12)



Strategy	Validation statistic		Estimate (95% CI) of mean $\mu$	95% Prediction interval	$P$ (%)	$\hat{\tau}$ 95% CI
	Calibration-in-the-large	Calibration-in-the-small				
<b>Strategy (1): Develop using logistic regression and implement with intercept estimated in external validation study</b>						
Calibration slope	-0.130 (-0.185, -0.075)	-0.195, -0.065	1	0.008	57	0.138
Log-expected/observed	0.095 (0.085, 0.197)	0.047, 0.128	0	0.0009	34	0.0017
C statistic	0.687 (0.670, 0.704)	0.645, 0.729	34	0.532	97	0.532
Calibration-in-the-large	-0.004 (-0.115, 0.305)	-1.240, 1.232	97	0.165	97	0.391
Calibration slope	0.980 (0.893, 1.077)	0.885, 1.375	59	0.087	640	0.734
Log-expected/observed	0.022 (-0.056, 0.250)	-0.485, 0.427	38	0.019	640	0.734
<b>Strategy (2): Develop using logistic regression and implement with average study intercept taken from developed model</b>						
Calibration slope	-0.047 (-0.120, 0.214)	-0.584, 0.678	89	0.270	59	0.167
Log-expected/observed	0.976 (0.851, 1.102)	0.578, 1.375	89	0.195	640	0.734
C statistic	-0.029 (-0.150, 0.093)	0.687 (0.669, 0.705)	38	0.019	640	0.734
<b>Strategy (3): Develop using logistic regression and implement with intercept taken from a study used in development data with a similar prevalence</b>						
Calibration slope	0.047 (-0.120, 0.214)	-0.584, 0.678	89	0.270	59	0.167
Log-expected/observed	0.976 (0.851, 1.102)	0.578, 1.375	89	0.195	640	0.734
C statistic	-0.029 (-0.150, 0.093)	0.687 (0.669, 0.705)	38	0.019	640	0.734

\* A bivariate meta-analysis was fitted to calibration-in-the-large calibration slope, and C statistic, and then again for (log(expected/observed)/calibration slope) and C statistic, respectively. Results were practically the same for calibration slope and C statistic, regardless of the bivariate model fitted.

## Example 2

### Prognosis of amyotrophic lateral disease

- IPD-MA
  - 14 cohort studies (specialized ALS centres)
- Sample size
  - 190 to 1,936 per study (total N = 11,475)
- Composite endpoint
  - Non-invasive ventilation for more than 23h/day, or death
  - Total number of events E = 8,819
- Median follow-up: 97.5 months

Development of the NCALS model



### Example 1

#### Diagnosis of deep vein thrombosis (N=12)

Table 2: Joint predicted probability of "good" discrimination and calibration performance of the DVT model for each of the three implementation strategies, derived using the multivariate meta-analysis results for the C statistic and calibration slope shown in Table 1						
Joint predicted probability of meeting criteria in new population						
Strategy (1): Development using logistic regression and implement with intercept estimated in external validation study		Strategy (2): Development using logistic regression and implement with average study intercept taken from developed model		Strategy (3): Development using logistic regression and implement with intercept taken from a study used in development data with a similar prevalence		
Calibration	Minimum C required	Calibration	Minimum C required	Calibration	Minimum C required	Calibration
slope required	with intercept estimated in external validation study	slope required	with intercept estimated in external validation study	slope required	with intercept estimated in external validation study	slope required
0.9-1.1	0.70	0.9-1.2	0.70	0.9-1.2	0.70	0.9-1.2
0.8-1.2	0.65	0.8-1.2	0.65	0.8-1.2	0.65	0.8-1.2
0.9-1.1	0.65	0.9-1.1	0.65	0.9-1.1	0.65	0.9-1.1
0.8-1.2	0.65	0.8-1.2	0.65	0.8-1.2	0.65	0.8-1.2

Abbreviation: DVT, deep vein thrombosis.

Notice that using the average intercept term is not problematic when the main focus is on calibration slope and C-statistic



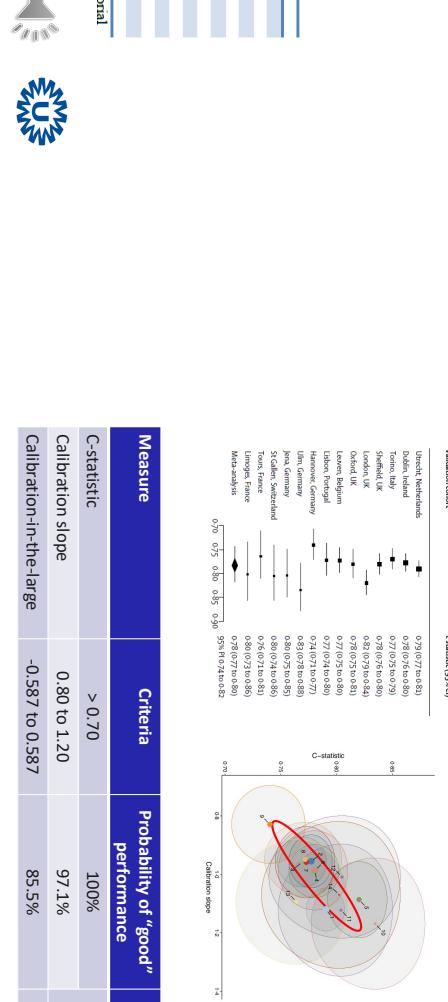
## Example 2

### Prognosis of amyotrophic lateral disease

- Royston-Parmar survival model with country-specific (but proportional) baseline hazard

Variable	Value
$\beta_0$	-5.409
$\beta_1$	2.645
$\beta_2$	-0.265
$\beta_3$	0.182
$\beta_4$ (ALSFRS-R slope)	-1.837
$\beta_5$ (Diagnostic delay)	-2.373
$\beta_6$ (Age at onset)	-0.267
$\beta_7$ (Forced vital capacity)	0.477
$\beta_8$ (Bulbar onset)	0.269
$\beta_9$ (Definite ALS*)	0.233
$\beta_{10}$ (Frontotemporal dementia)	0.388
$\beta_{11}$ (Corticosteroid repeat expansion)	0.256

Supplementary Table S15. Parameters of the final prediction model. \*According to the El Escorial criteria.



## Example 2

### Example 2

### THE LANCET Neurology

Volume 17, Issue 5, May 2018, Pages 423-433

Articles

- Prognosis for patients with amyotrophic lateral sclerosis: development and validation of a personalised prediction model

Henk-Jan Westenberg MD<sup>a</sup>, Thomas P A Debrij PhD<sup>a, b</sup>, Anne E Visser MD<sup>a</sup>, Ruben P A van Elst MD<sup>a</sup>, James P K Roosje MSc<sup>a</sup>, Andrea Calvo MD<sup>a</sup>, Sarah Larin BSc<sup>c</sup>, Prof Christopher J McDermott PhD<sup>b</sup>, Alexander G Thompson BM&ChB<sup>b</sup>, Susana Pinto PhD<sup>b</sup>, Xanthe Kolevaeva MD<sup>b</sup>, Angela Rosenblum MD<sup>b</sup>, Beatrice Stuhendorff PhD<sup>b</sup>, Helma Sommer<sup>b</sup>, Bas M Middelkoop<sup>b</sup>, Anneke M Dekker MD<sup>a</sup>, Joke J F A van Vugt PhD<sup>a</sup>, Wouter van Rijen MD<sup>a</sup> ... Prof Leonard H van den Berg MD<sup>a, b, c</sup>

### Internal-external cross-validation

### Example 2



## Example 2

The life expectancy of Stephen Hawking, according to the ENCALS model

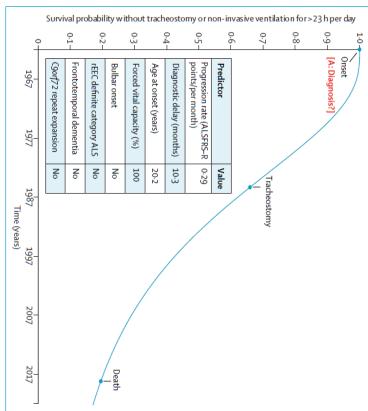


Figure: Personalised survival curve for Stephen Hawking  
Using publicly available data of disease characteristics at diagnosis in 1963, we were able to estimate the probability to survive without tracheostomy or non-invasive ventilation for more than 23 hours per day. The predictor values used for calculating the survival probability are summarised in the table.  
ALS=amyotrophic lateral sclerosis; ALSFS=ALS functional rating scale; REF=repeat Expansion criteria.



## Example 2

The life expectancy of Stephen Hawking, according to the ENCALS model

"Using publicly available data, we examined whether Professor Hawking's survival was as rare as his intellectual performance, or could be predicted solely based on his disease characteristics at diagnosis in 1963."

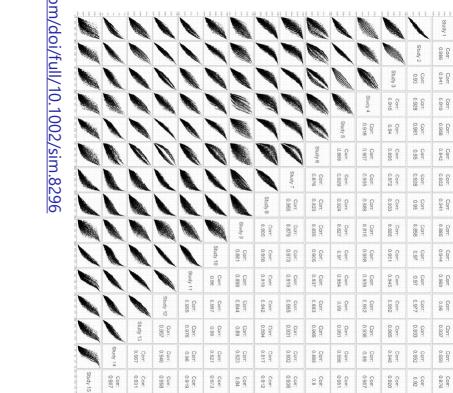
- Predicted 10-year survival probability: 94%
- The IQR for his predicted survival lay between 1981 and 2011
- Young age of onset was the most important factor for his long survival

## Assessing heterogeneity in predictor effects

### Assessing heterogeneity in predictor effects

- Generalizable prediction models have little heterogeneity in predictor effects
- However, heterogeneity may sometimes appear because of collinearity
- Heterogeneity in absolute risk is what matters most
- This can be explored by plotting the predictions of study-specific models in a pairwise comparison

<https://onlinelibrary.wiley.com/doi/10.1002/sim.8296>



## Summary points

Major advantages of large clustered datasets

- Improve the performance of novel prediction models across different study populations
  - Attain a better understanding of the generalizability of a prediction model
  - Explore heterogeneity in model performance and the added value of a novel (bio)marker

Unfortunately, most researchers analyze their IPD as if representing a single dataset!



## Remaining challenges in IPD-MA

- IPD-MA not a solution for poorly designed primary studies
  - Prospective multi-center studies remain important
- Synthesis strategies from intervention research cannot directly be applied in prediction research (due to focus on absolute risks)
  - Adjustment to local circumstances often needed
  - One model fits all?
  - Methods for tailoring still underdeveloped

New methods are on their way!

## Assessing heterogeneity in predictor effects

Received: 4 June 2017 | Revised: 12 March 2019 | Accepted: 4 June 2019

DOI: 10.1002/sim.8296

Diagnostic and Prognostic Research

https://doi.org/10.1186/s41512-019-0059-4

(2019) 3:13

Open Access

CrossMark

Check for updates

## Summary points

Debray et al. Diagnostic and Prognostic Research

(2019) 3:13

Diagnostic and Prognostic Research

METHODOLOGY

## Evidence synthesis in prognosis research

Thomas P.A. Debray<sup>1,2\*</sup>, Valentijn M.T. de Jong<sup>1†</sup>, Karel G.M. Moons<sup>1,2</sup> and Richard D. Riley<sup>3</sup>

### Abstract

Over the past few years, evidence synthesis has become essential to investigate and improve the generalizability of medical research findings. This strategy often involves a meta-analysis to formally summarize quantities of interest, such as relative treatment effect estimates. The use of meta-analysis methods is, however, less straightforward in prognosis research because substantial variation exists in research objectives, analysis methods and the level of reported evidence.

We present a generic overview of statistical methods that can be used to summarize data of prognostic factor and prognostic model studies. We discuss how aggregate data of individual participant data, or a combination thereof can be combined through meta-analysis methods. Recent examples are provided throughout to illustrate the various methods.

**Keywords:** Prediction, Meta-analysis, Prognosis, Validation, IPD

<https://onlinelibrary.wiley.com/doi/full/10.1002/sim.8296>



## R software

[metanisc: Diagnostic and Prognostic Meta-Analysis](#)

Meta-analysis of diagnostic and prognostic modeling studies. Summarize estimates of prognostic factors, diagnostic test accuracy and prediction model performance. Validate, update and combine published prediction models. Develop new prediction models with data from multiple studies.

Version:

0.1.9

Depends:

R (>= 2.0), stats, graphics

Imports:

[metabias](#) (>= 0.0), [minmeta](#), [ellipses](#), [lme4](#), [nlme](#), [sampleSize](#)

Suggests:

[lumix](#), [RMS](#), [textual](#) (>= 1.0.2)

Published:

2018-05-13

Author:

Thomas Debray  [aut, cre], Valentin de Jong [au]

Maintainer:

Thomas Debray <thomas.debray@gmail.com>

License:

[GPL-3](#)

NeedsCompilation:

no

In views:

MetaAnalysis

CRAN checks:

[metanisc results](#)

Downloads :

Reference manual: [metanisc.pdf](#)

Package source: [metanisc\\_0.1.9.tar.gz](#)

Windows binary: [r-devel](#); [metanisc\\_0.1.9.zip](#); [r-release](#); [metanisc\\_0.1.9.zip](#)

OS X binaries: [r-release](#); [metanisc\\_0.1.9.tgz](#); [r-devel](#); [metanisc\\_0.1.9.zip](#)

Old sources: [metanisc archive](#)

Linking:

Please use the canonical form <https://CRAN.R-project.org/package=metanisc> to link to this page.

