

Better predictions using big(ger) datasets

Thomas Debray, PhD
Assistant Professor



Risk prediction

- Risk prediction = foreseeing / foretelling
... (probability) of something that is yet unknown
- Turn available information (predictors) into a statement about the probability:
 - ... of having a particular disease -> **diagnosis**
 - ... of developing a particular event -> **prognosis**



Why do we predict?

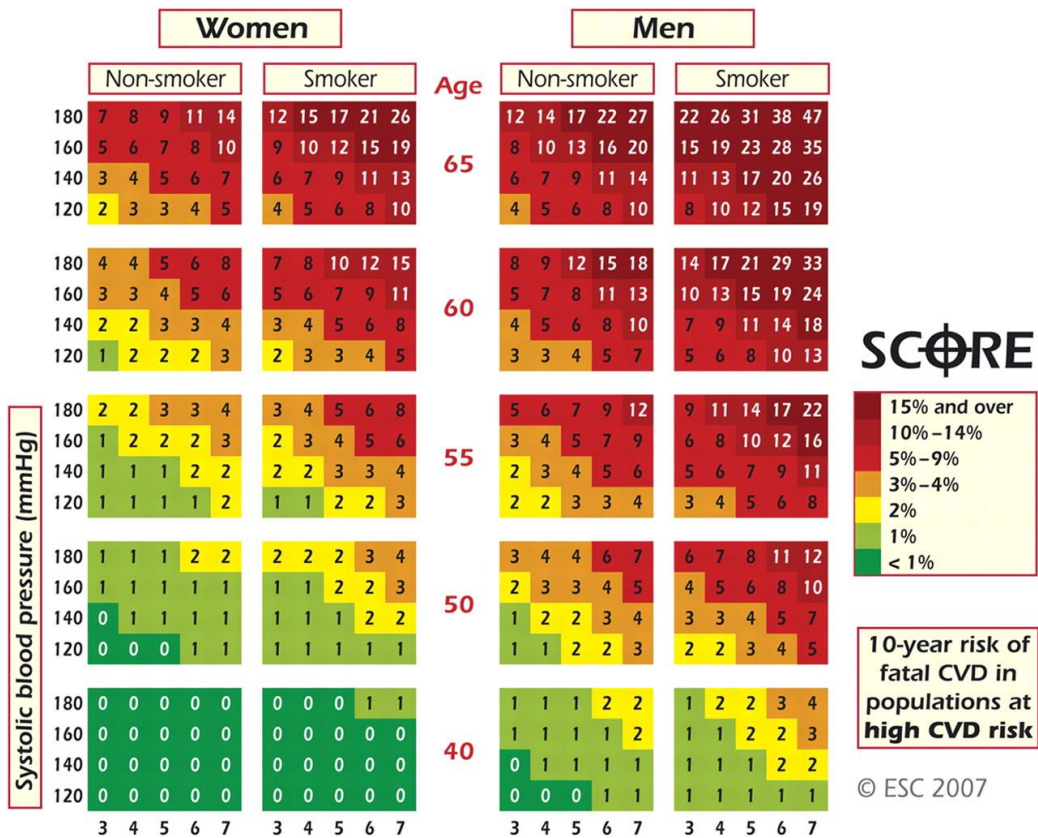
- Identification of high risk individuals
 - To inform patients and their families
 - To guide treatment decisions (“**precision medicine**”)
 - To design randomized trials
- Data analysis
 - To deal with missing values
 - To match/subclassifiy patients
 - ...



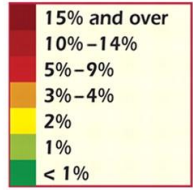
How do we predict?

- Combine information from multiple predictors
 - Patient characteristics (e.g. age, gender)
 - History and physical examination results (e.g. blood pressure)
 - Imaging results
 - (Bio)markers (e.g. coronary plaque)
- Develop a multivariable statistical model
 - Need for individual participant data (e.g. from cohort studies)
 - Many strategies available (e.g. logistic regression)





SCORE



10-year risk of fatal CVD in populations at high CVD risk

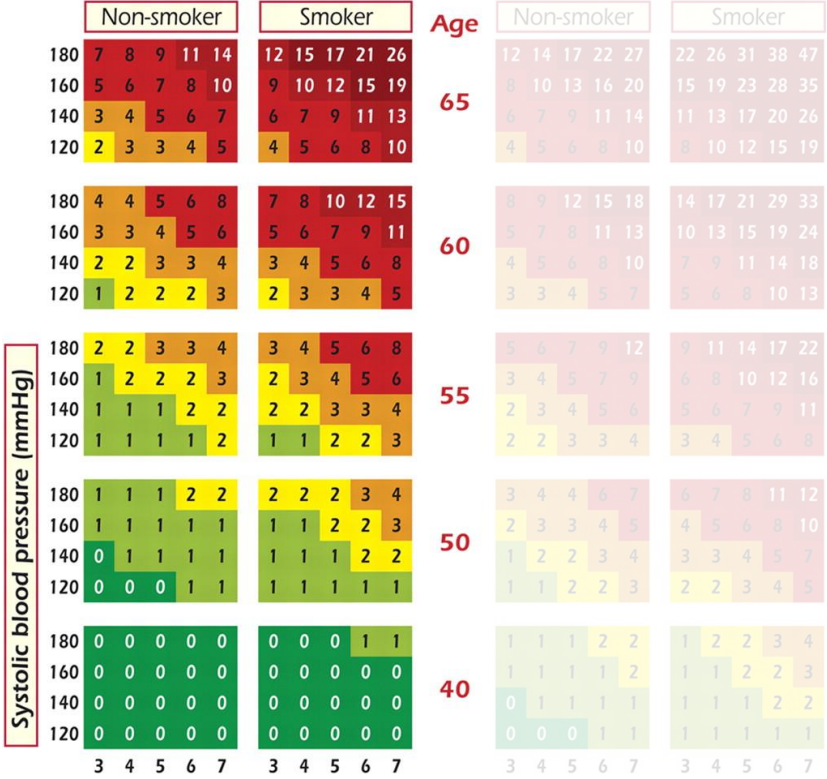
© ESC 2007

Total cholesterol: HDL Cholesterol ratio

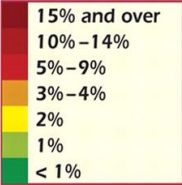


Women

Men



SCORE



10-year risk of fatal CVD in populations at high CVD risk

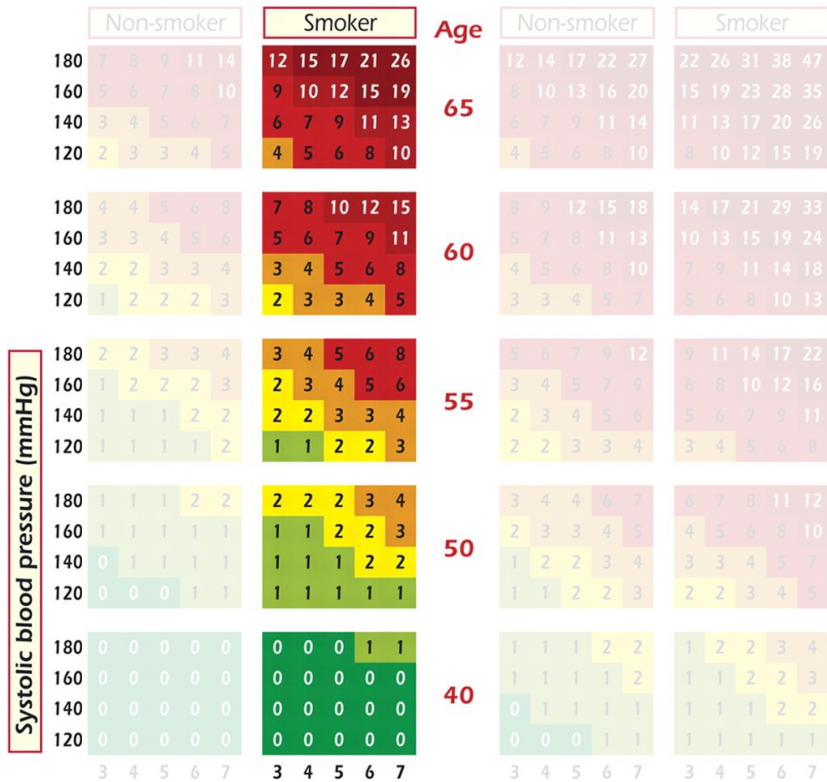
© ESC 2007

Total cholesterol: HDL Cholesterol ratio

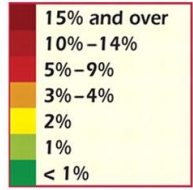


Women

Men



SCORE



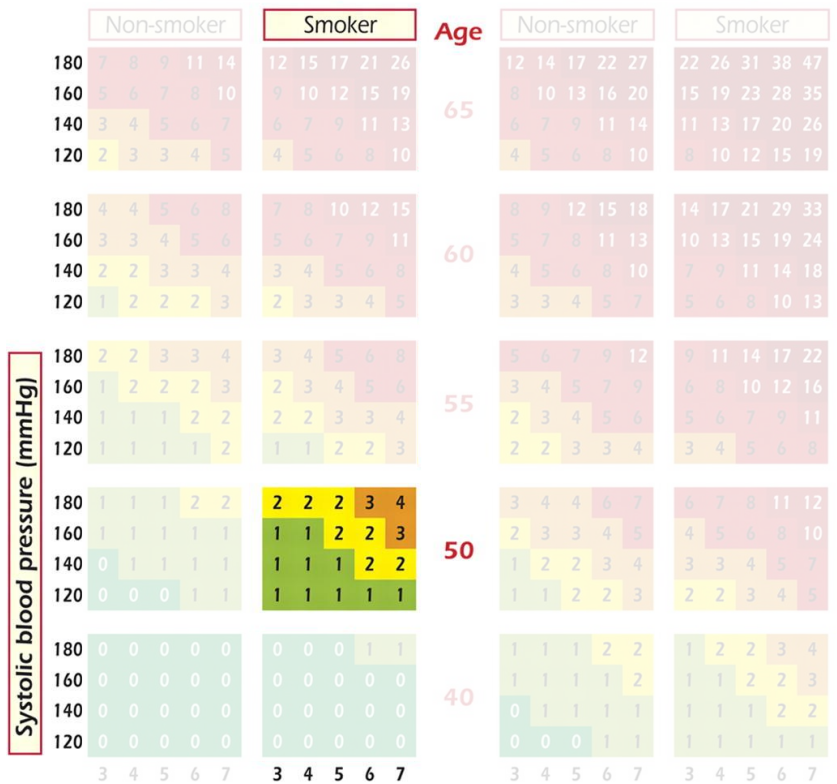
10-year risk of fatal CVD in populations at high CVD risk

© ESC 2007

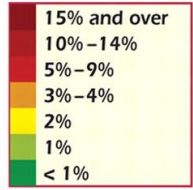


Women

Men



SCORE



10-year risk of fatal CVD in populations at high CVD risk

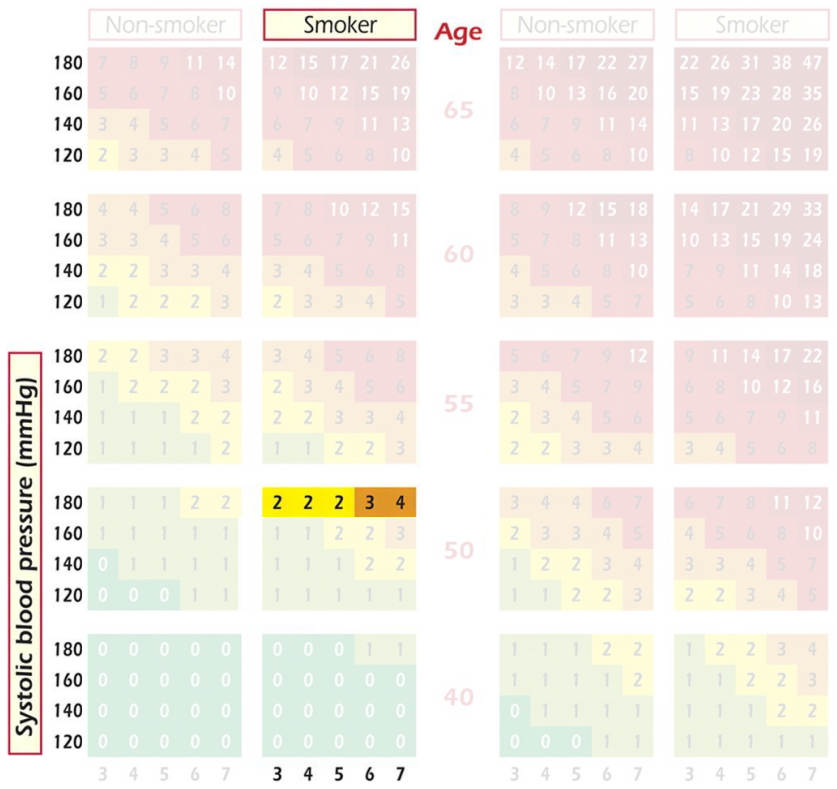
© ESC 2007

Total cholesterol: HDL
Cholesterol ratio

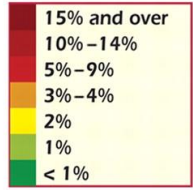


Women

Men



SCORE



10-year risk of fatal CVD in populations at high CVD risk

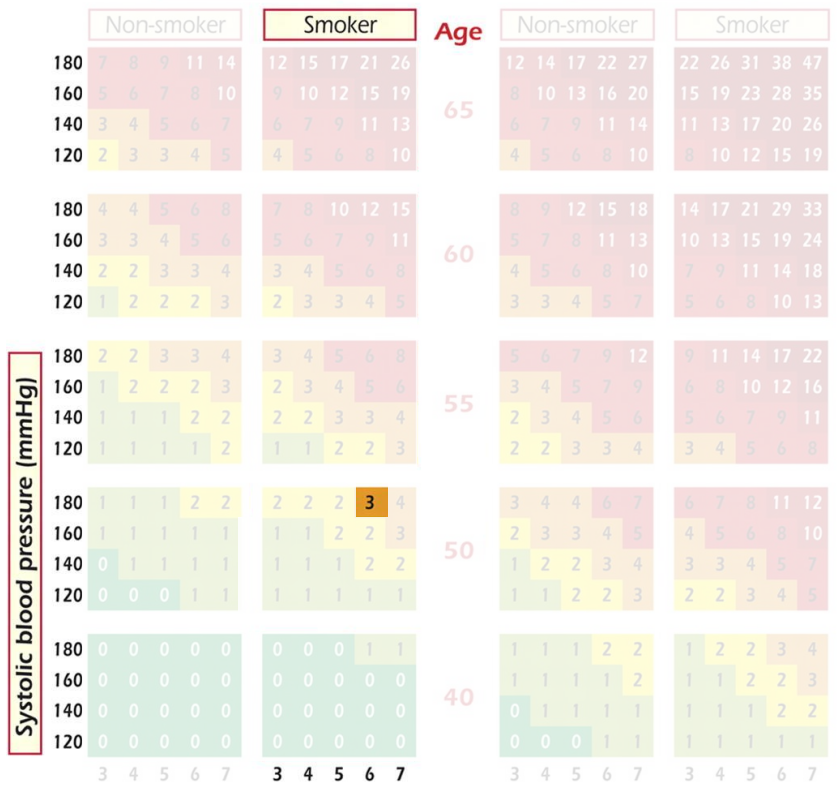
© ESC 2007

Total cholesterol: HDL
Cholesterol ratio

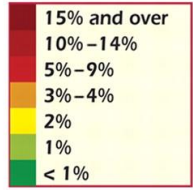


Women

Men



SCORE



10-year risk of fatal CVD in populations at high CVD risk

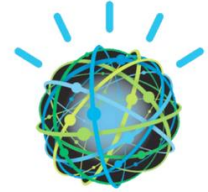
© ESC 2007

Total cholesterol: HDL Cholesterol ratio



Watson for Oncology

“Bring personalized, evidence-supported cancer care plans to your patients”



IBM WATSON

- Interpret cancer patients' clinical information
- Digest doctor's notes, medical studies, and clinical guidelines
- Provide individualized treatment recommendations
- Adopted by more than 150 hospitals and healthcare organizations

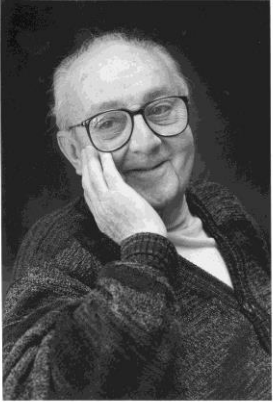


Hype meets reality

- Focus on US clinical practice and demographics
- Reliance on varies among hospitals
- Multiple examples of unsafe and incorrect treatment recommendations
- Lack of validation by independent scientists
- Lack of clinical trials to assess effectiveness



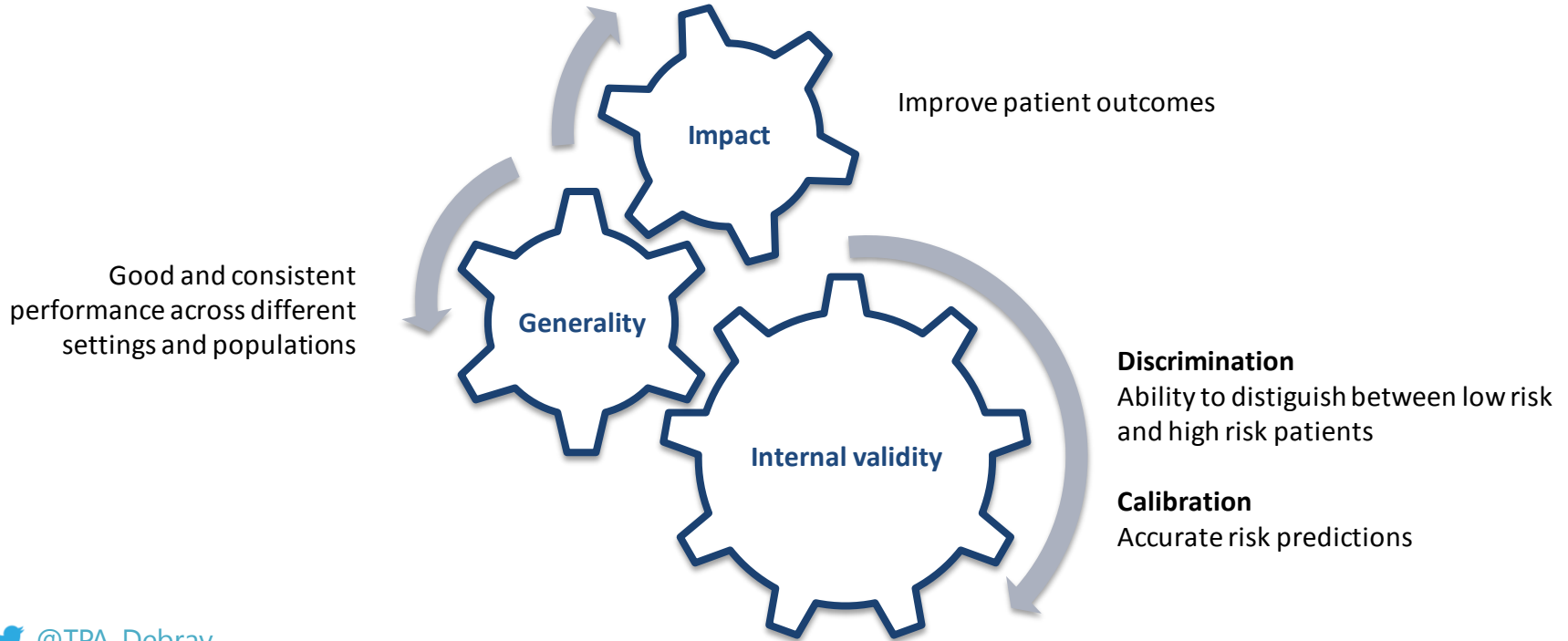
Most models are not as good as we think



“All models are wrong, but some are useful”

George Box

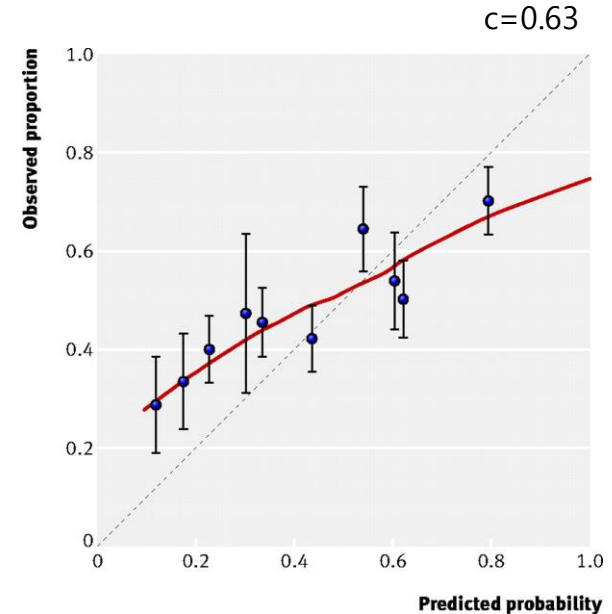
What is a “good” prediction model?



What is a “good” prediction model?

Common measures of prediction model performance

- Discrimination
 - Concordance (c-) statistic
 - Area under the ROC curve
- Calibration
 - Calibration intercept (calibration-in-the-large)
 - Calibration slope
 - Ratio of expected versus observed events



Most models are not as good as we think

How to assess and improve the generalizability of prediction models?



The rise of “big” data sets



The rise of “big” data sets

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

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





Validation of existing prediction models

- Assess prediction model performance **multiple times**
 - In new patients from the same (target) population
 - In new patients from different (but related) populations
 - To evaluate generalizability across different settings and populations
- Meta-analysis methods are needed
 - To summarize prediction model performance
 - To investigate sources of between-study heterogeneity



Validation of QRISK 2

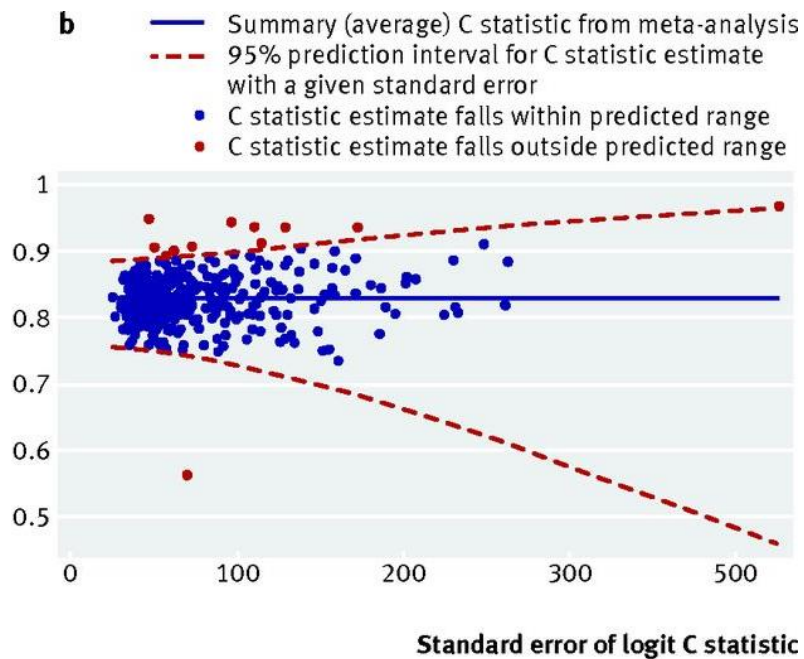
Registry data with 1.58 million patients from 364 practices

Objective To evaluate the performance of the QRISK2 score for predicting 10-year cardiovascular disease in an independent **UK cohort of patients from general practice records** and to compare it with the NICE version of the Framingham equation and QRISK1.

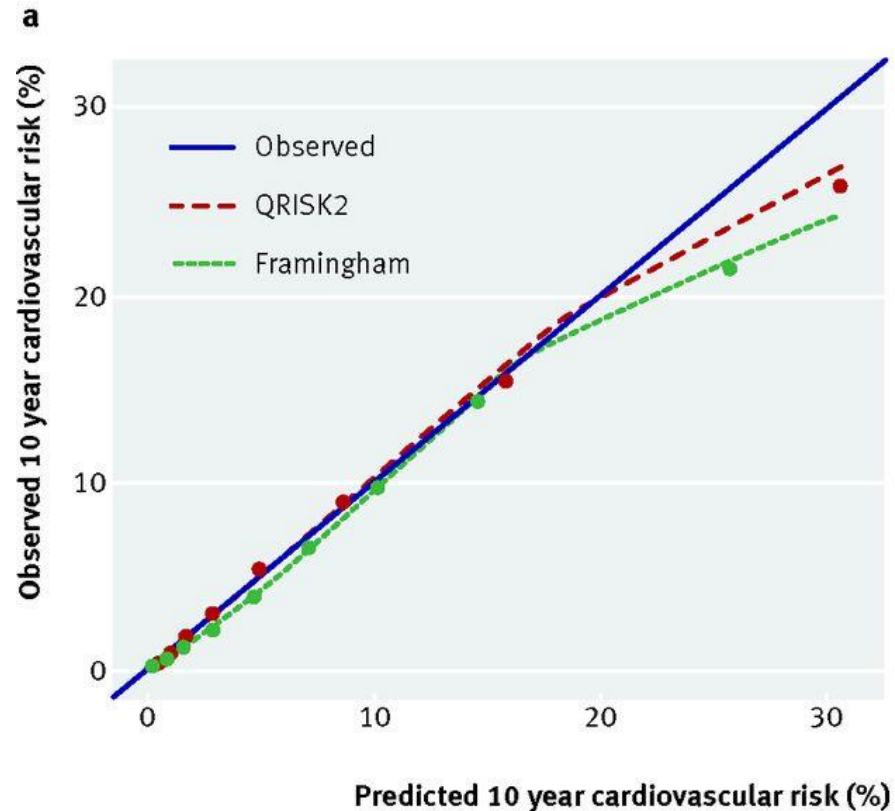
Design Prospective cohort study to validate a cardiovascular risk score.

Setting 365 practices from United Kingdom contributing to The Health Improvement Network (THIN) database.

Participants **1.58 million patients** registered with a general practice between 1 January 1993 and 20 June 2008, aged 35-74 years (9.4 million person years) with 71 465 cardiovascular events.



Summary (average) C statistic = 0.83 (95% CI 0.826 to 0.833)
 95% prediction interval for true C statistic in a new practice = 0.76 to 0.88



Key references

- Debray et al. A framework for meta-analysis of prediction model studies with binary and time-to-event outcomes. [Stat Methods Med Res 2018](#)
- Debray et al. A guide to systematic review and meta-analysis of prediction model performance. [BMJ 2017](#)
- 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](#)
- Snell et al. Multivariate meta-analysis of individual participant data helped externally validate the performance and implementation of a prediction model. [J Clin Epidemiol 2015](#)



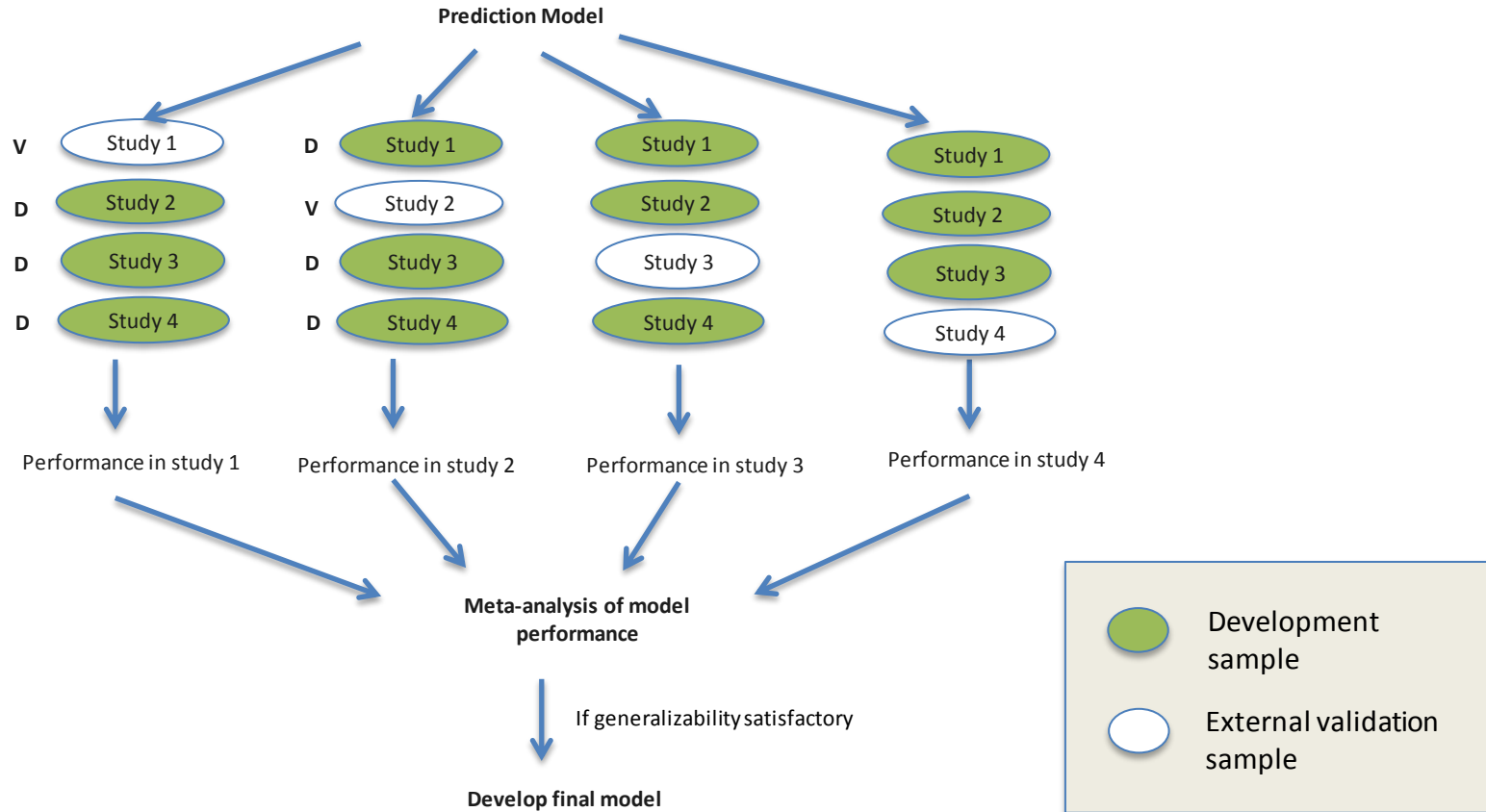


Development of prediction models

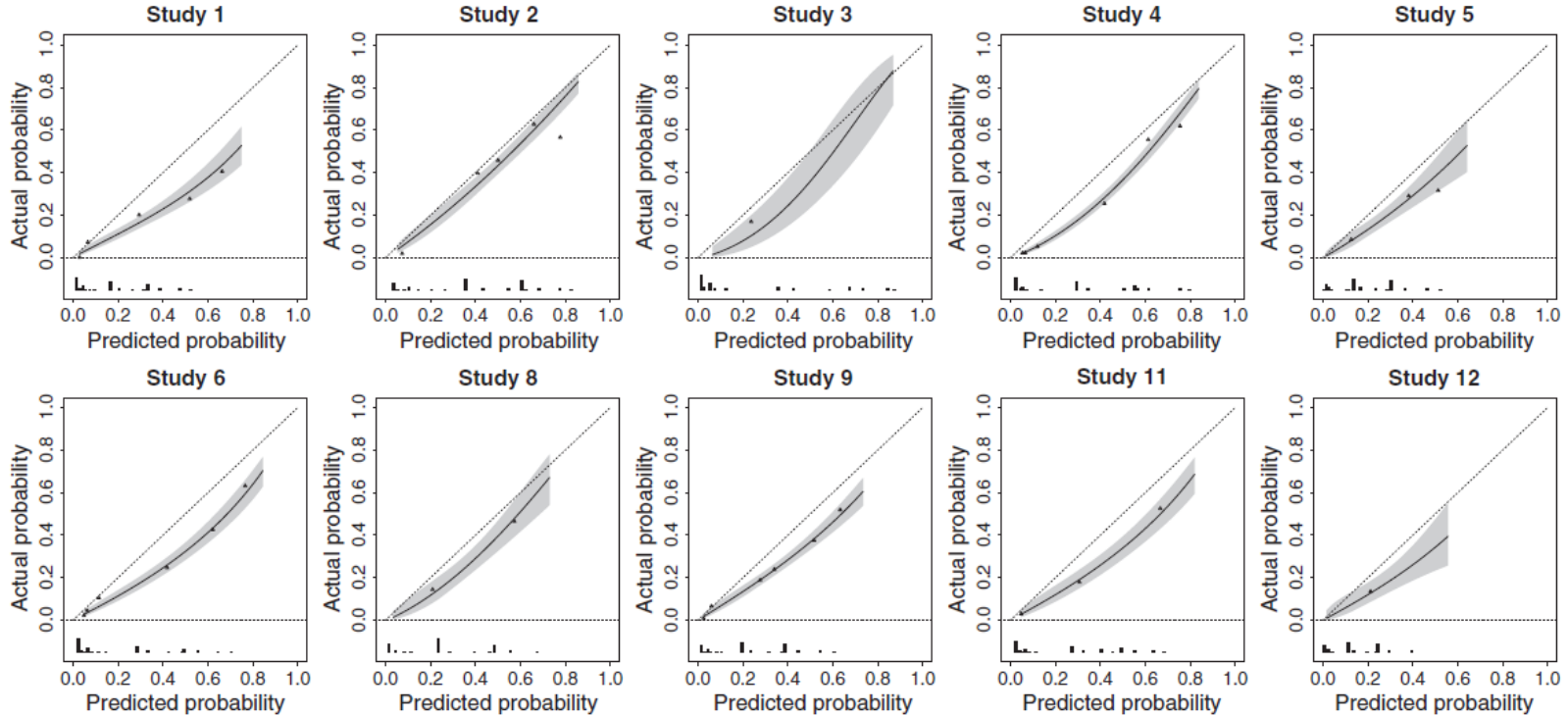
- Access to big(ger) datasets enables to assess model transportability (rather than merely reproducibility) during its development
- Identify and account for heterogeneity in prediction model performance via internal-external cross-validation
- This allows to optimize prediction model generalizability



Internal-external cross-validation (IECV)



IECV allows for many external validations

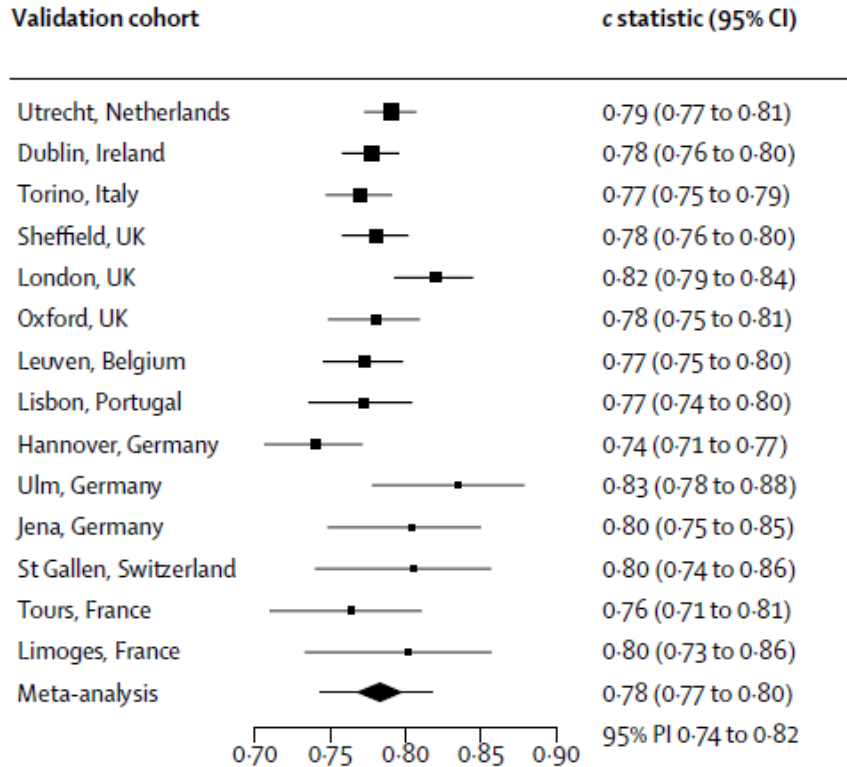


Development and validation of ENCALs

Prognosis for patients with amyotrophic lateral sclerosis (ALS)

- Cohort data from 11,475 patients from 14 European ALS centres
- Composite survival outcome (non-invasive ventilation for more than 23 h per day, tracheostomy, or death)
- Development of multivariable Royston-Parmar models
- Assessment of generalizability via IECV

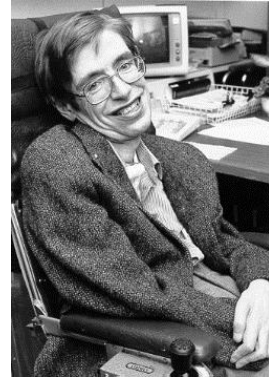
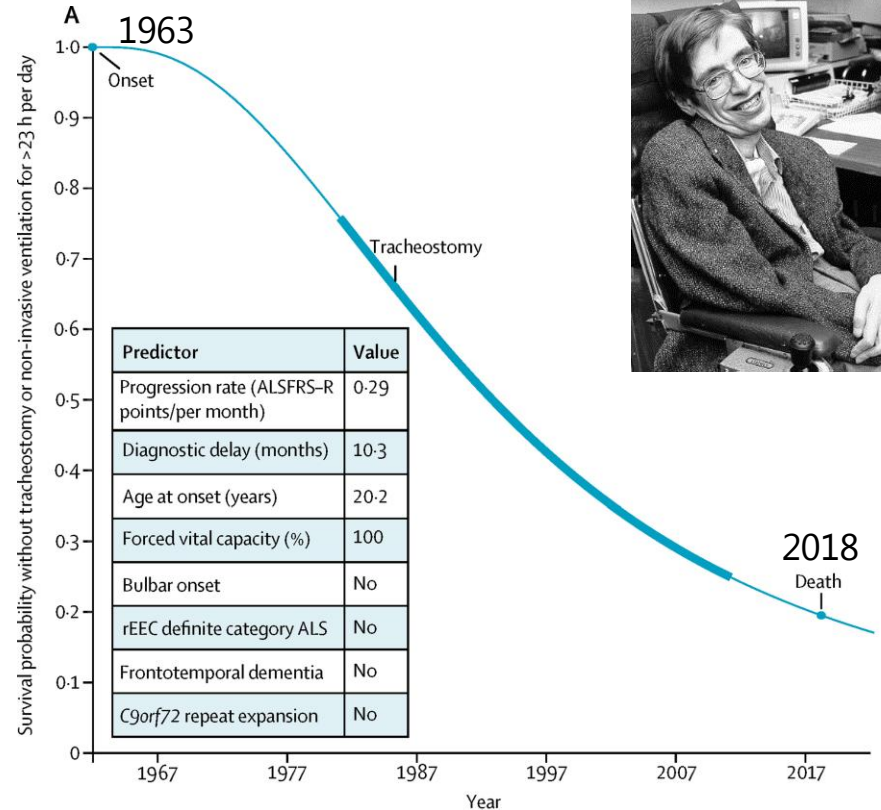
Development and validation of ENCALs



Pr(c>0.70)=100%

The life expectancy of Stephen Hawking

- Prediction in 1963: 2 year survival
- Prediction according to ENCALs: **94%** to survive at least 10 years
- Young age of onset was the most important factor for his long survival



Key references


- Ahmed et al. Developing and validating risk prediction models in an individual participant data meta-analysis. [BMC Med Res Methodol 2014](#)
- Debray et al. A framework for developing, implementing, and evaluating clinical prediction models in an individual participant data meta-analysis. [Stat Med 2013](#)
- Debray et al. Individual Participant Data (IPD) Meta- analyses of Diagnostic and Prognostic Modeling Studies: Guidance on Their Use. [PLoS Med 2015](#)

Future developments

Software

metamisc: 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 ($\geq 3.2.0$), stats, graphics
Imports: [metafor](#) ($\geq 2.0.0$), [mvtnorm](#), [ellipse](#), [lme4](#), [plyr](#), [ggplot2](#)
Suggests: [runjags](#), [rjags](#), [testthat](#) ($\geq 1.0.2$)
Published: 2018-05-13
Author: Thomas Debray  [aut, cre], Valentijn de Jong [aut]
Maintainer: Thomas Debray <thomas.debray at gmail.com>
License: [GPL-3](#)
URL: <http://r-forge.r-project.org/projects/metamisc/>
NeedsCompilation: no
In views: [MetaAnalysis](#)
CRAN checks: [metamisc results](#)

Downloads:

Reference manual: [metamisc.pdf](#)
Package source: [metamisc 0.1.9.tar.gz](#)
Windows binaries: r-devel: [metamisc 0.1.9.zip](#), r-release: [metamisc 0.1.9.zip](#), r-oldrel: [metamisc 0.1.9.zip](#)
OS X binaries: r-release: [metamisc 0.1.9.tgz](#), r-oldrel: [metamisc 0.1.8.tgz](#)
Old sources: [metamisc archive](#)

Linking:

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



Guidance and methods

- **Prognostic Research in Health Care: concepts, methods and impact**
editors: Richard Riley, Danielle Van der Windt, Peter Croft, Karel Moons
- **Evidence synthesis using individual participant data: Concepts, Methods and Guidance for Clinical Research**
editors: Richard Riley, Jayne Tierney, Lesley Stewart
- **Handbook of Meta-analysis**
editors: Christopher Schmid, Theo Stijnen, Ian White