



UMC Utrecht  
**Julius Center**

# Predicting modeling using large clustered data sets

Thomas Debray, PhD

Julius Center for Health Sciences and Primary Care  
Utrecht University & Cochrane Netherlands  
Utrecht, The Netherlands



# Contents

Introduction

Common pitfalls

Opportunities

- Evidence synthesis
- Big data
- Machine Learning



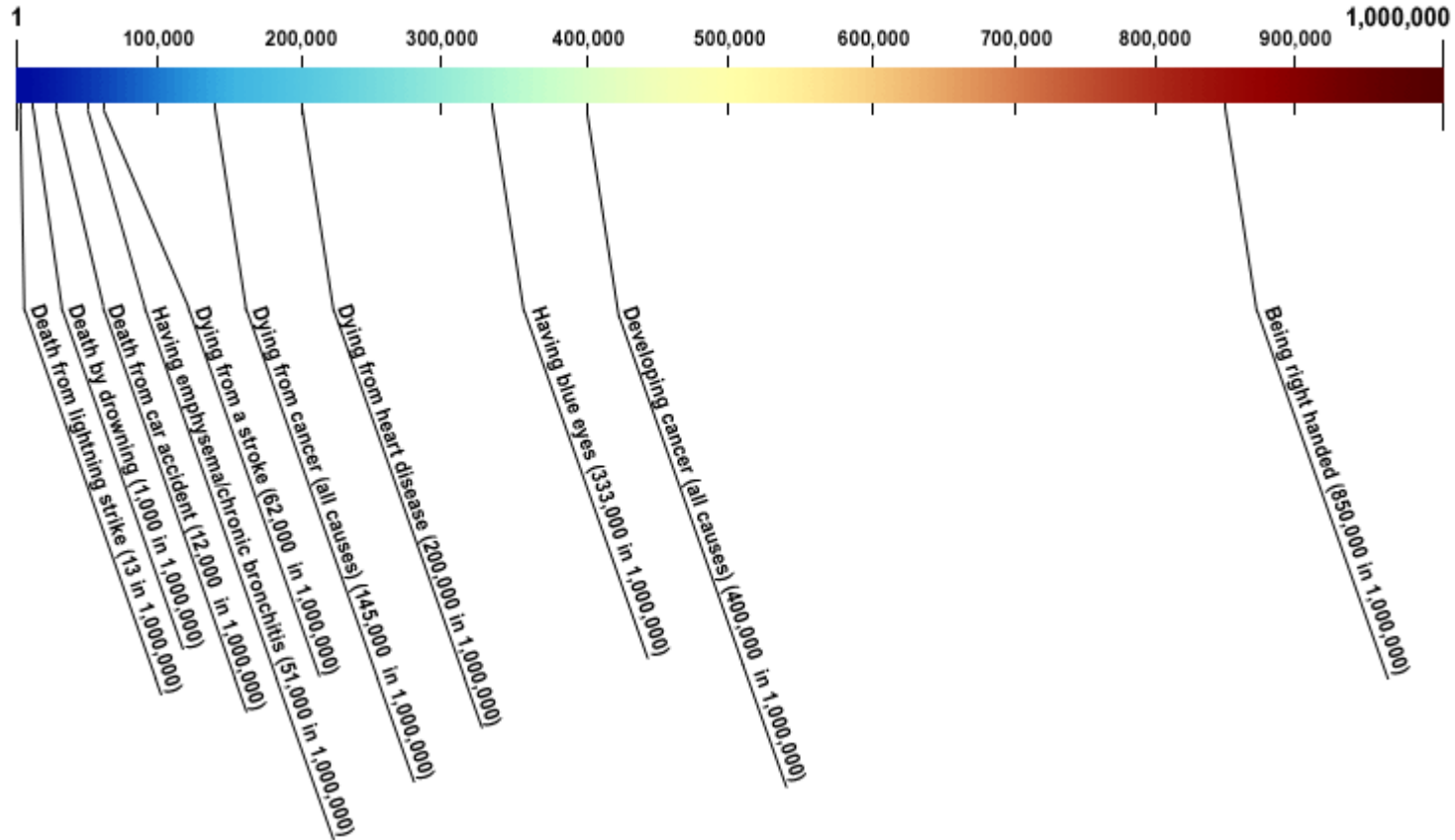
UMC Utrecht

# Prediction

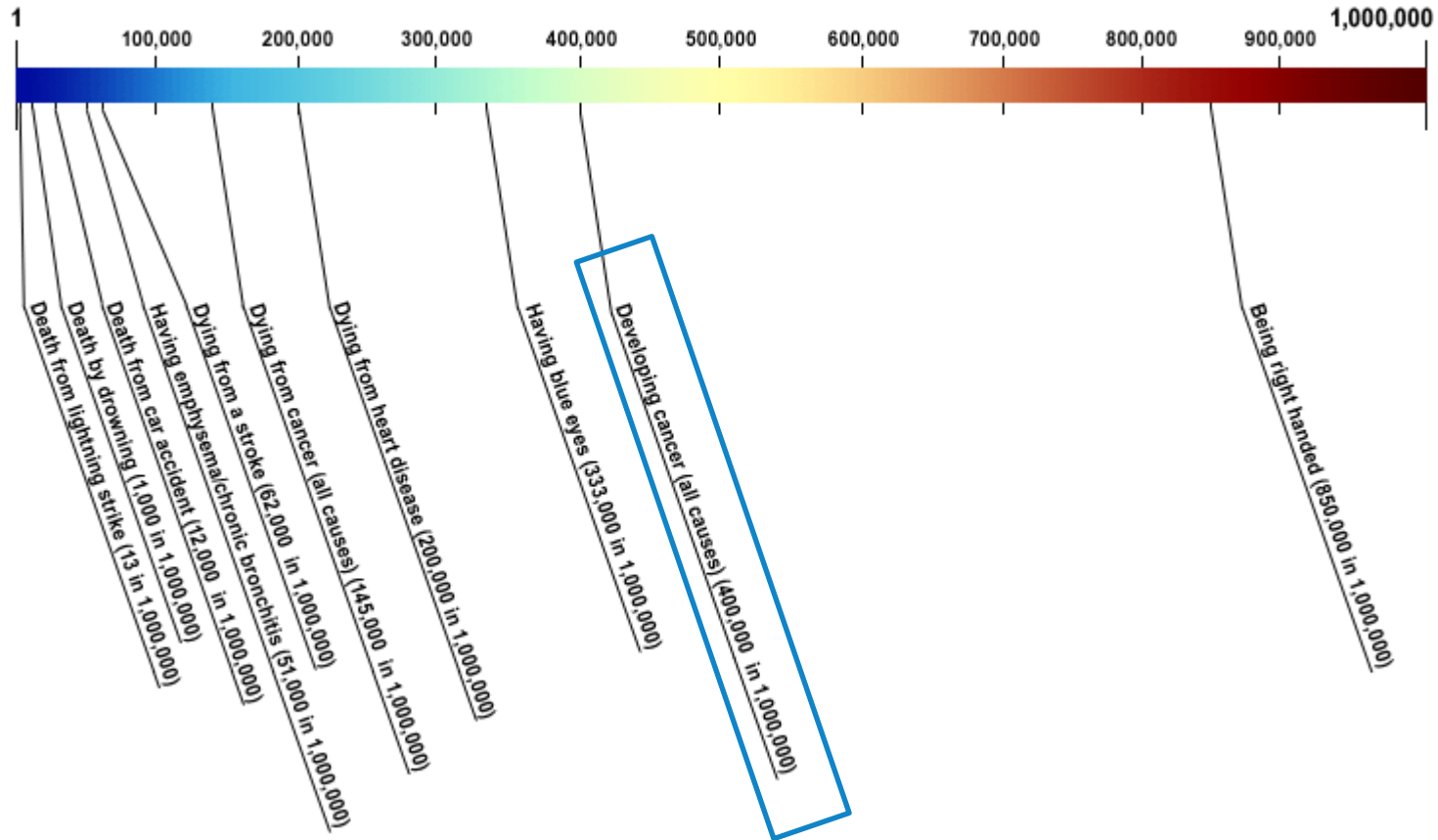
- Risk prediction = foreseeing / foretelling  
... (probability) of something that is yet unknown
- Turn available information (predictors) into a statement about the probability:
  - ... of having a particular disease -> **diagnosis**
  - ... of developing a particular event -> **prognosis**
- Use of prognostic information:
  - to inform patients and their families
  - to guide treatment and other clinical decisions
  - to create risk groups
  - ...



# Statistical Probabilities



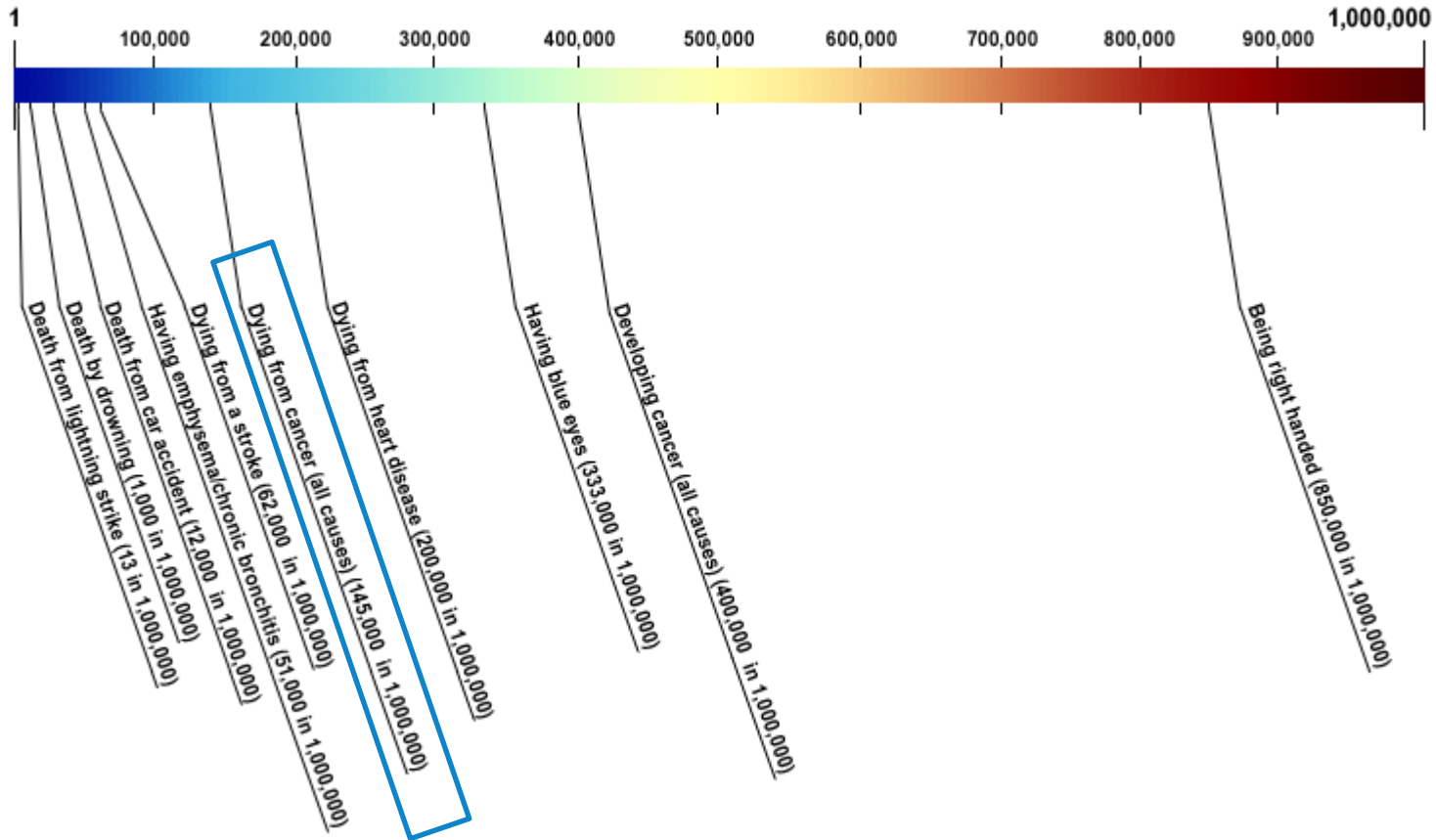
# Statistical Probabilities



Risk of developing cancer



# Statistical Probabilities

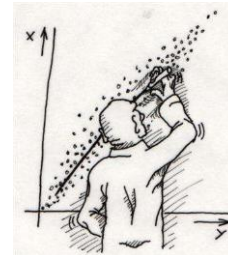
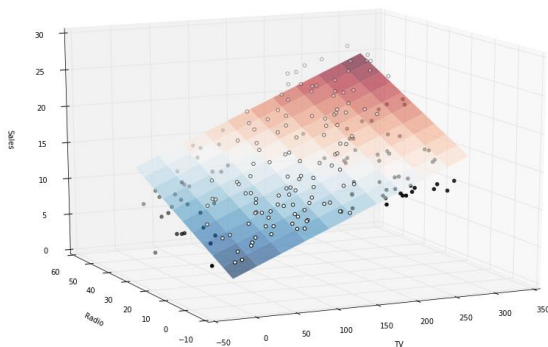


Risk of dying from cancer



# How do we predict?

- Combine information from multiple predictors
  - Subject 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 patient data from large cohort studies
  - Many strategies available (Regression, decision trees, neural networks, ...)





# Breast Cancer Risk Assessment Tool

An interactive tool to help estimate a woman's risk of developing breast cancer



Last modified date: 05/16/2011

> **Get Started with the Risk Tool**

About the Tool

Breast Cancer Risk Factors

Download Source Code

#### Page Options

 Print Page

#### Quick Links

[Breast Cancer Home Page](#)

[Breast Cancer: Prevention, Genetics, Causes](#)

[Current Clinical Trials: Breast Cancer In Situ: Treatment](#)

[Current Clinical Trials: Breast Cancer Prevention](#)

[Current Clinical Trials: Breast Cancer Screening](#)

[Breast Cancer Risk in American Women](#)



#### Need Help?

Contact us by phone, Web, and e-mail  
**1-800-4-CANCER**

The Breast Cancer Risk Assessment Tool is an interactive tool designed by scientists at the National Cancer Institute (NCI) and the [National Surgical Adjuvant Breast and Bowel Project \(NSABP\)](#) to estimate a woman's risk of developing [invasive breast cancer](#). See [About the Tool](#) for more information.

The Breast Cancer Risk Assessment Tool may be updated periodically as new data or research becomes available.

## Risk Tool

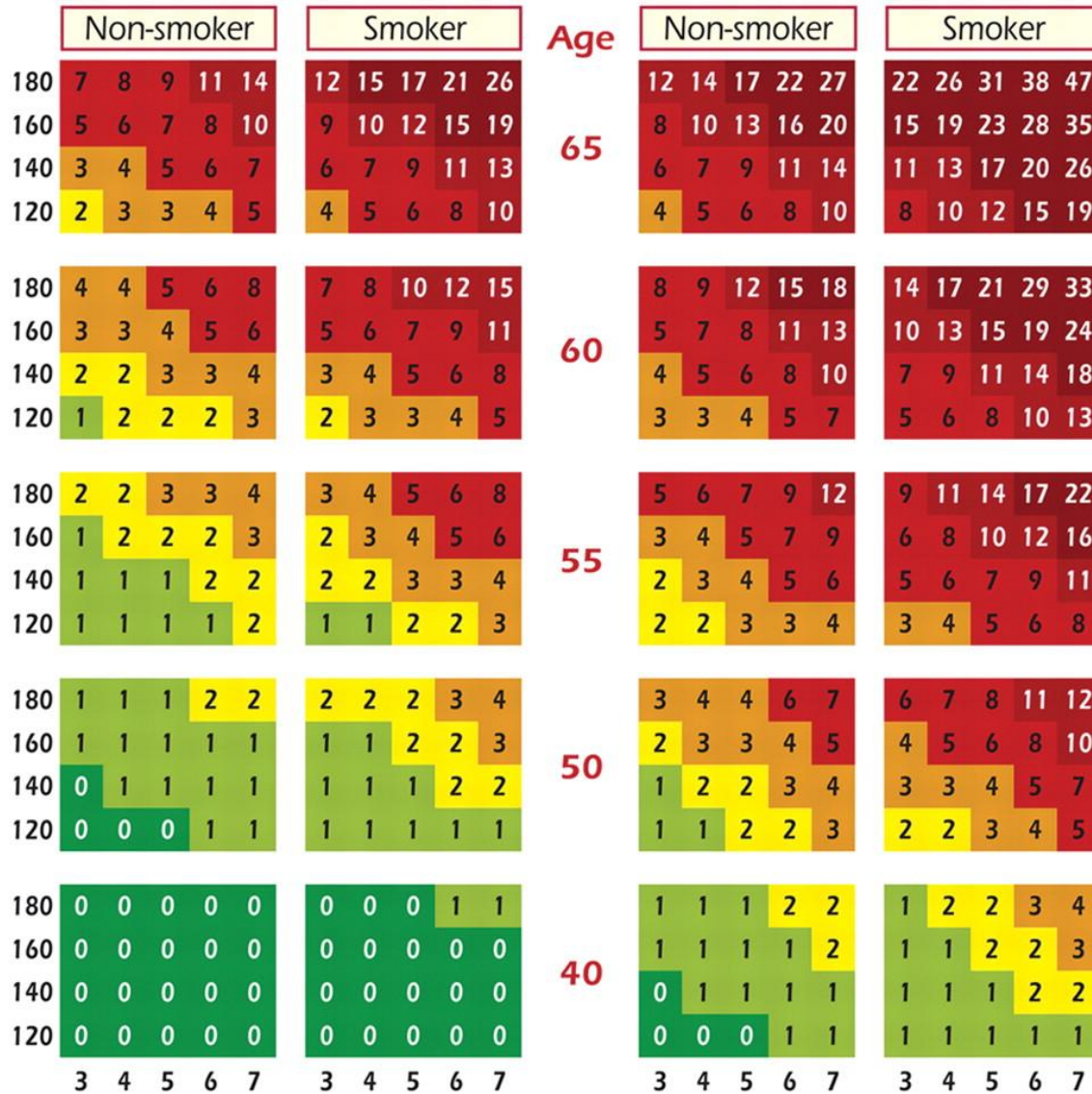
(Click a question number for a brief explanation, or [read all explanations](#).)

1. Does the woman have a medical history of any breast cancer or of ductal carcinoma in situ (DCIS) or lobular carcinoma in situ (LCIS) or has she received previous radiation therapy to the chest for treatment of Hodgkin lymphoma?
2. Does the woman have a mutation in either the *BRCA1* or *BRCA2* gene, or a diagnosis of a genetic syndrome that may be associated with elevated risk of breast cancer?
3. What is the woman's age?  
*This tool only calculates risk for women 35 years of age or older.*
4. What was the woman's age at the time of her first menstrual period?
5. What was the woman's age at the time of her first live birth of a child?
6. How many of the woman's first-degree relatives - mother, sisters, daughters - have had breast cancer?
7. Has the woman ever had a breast biopsy? 
  - 7a. How many breast biopsies (positive or negative) has the woman had?
  - 7b. Has the woman had at least one breast biopsy with atypical hyperplasia?
8. What is the woman's race/ethnicity? 
  - 8a. What is the sub race/ethnicity?

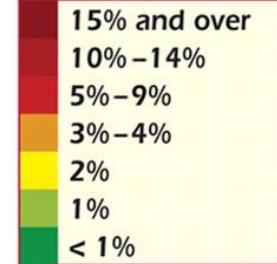
**Calculate Risk >**

## Women

## Men



# SCORE



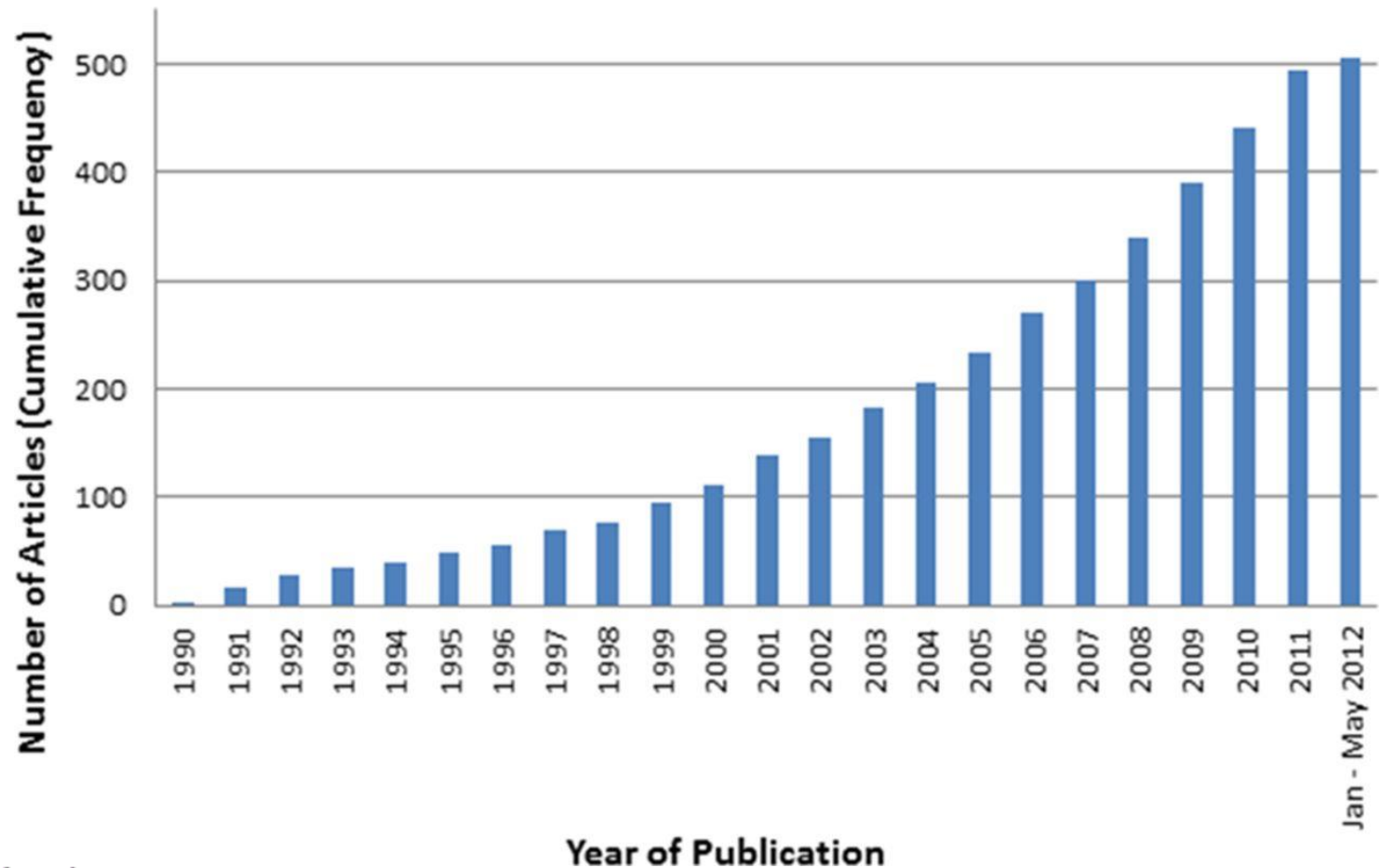
10-year risk of fatal CVD in populations at high CVD risk

© ESC 2007

Total cholesterol: HDL  
Cholesterol ratio

# Why focus on prediction models?

Cumulative growth in published CPM articles over time



# What is a good model?



# Phases of prediction model evaluation

Series in BMJ 2009 and in Heart 2012, Moons et al.

## Development

- Identify predictors
- Model building
- Internal validation

## Validation

- Performance in new individuals
- Narrow validation
- Broad validation

## Updating

- Adjust existing model to other settings/ populations to improve predictive performance

## Impact

- Quantify impact of use of model on decision making and health outcomes
- Experimental design

## Dissemination Implementation

- Widespread use
- Barriers

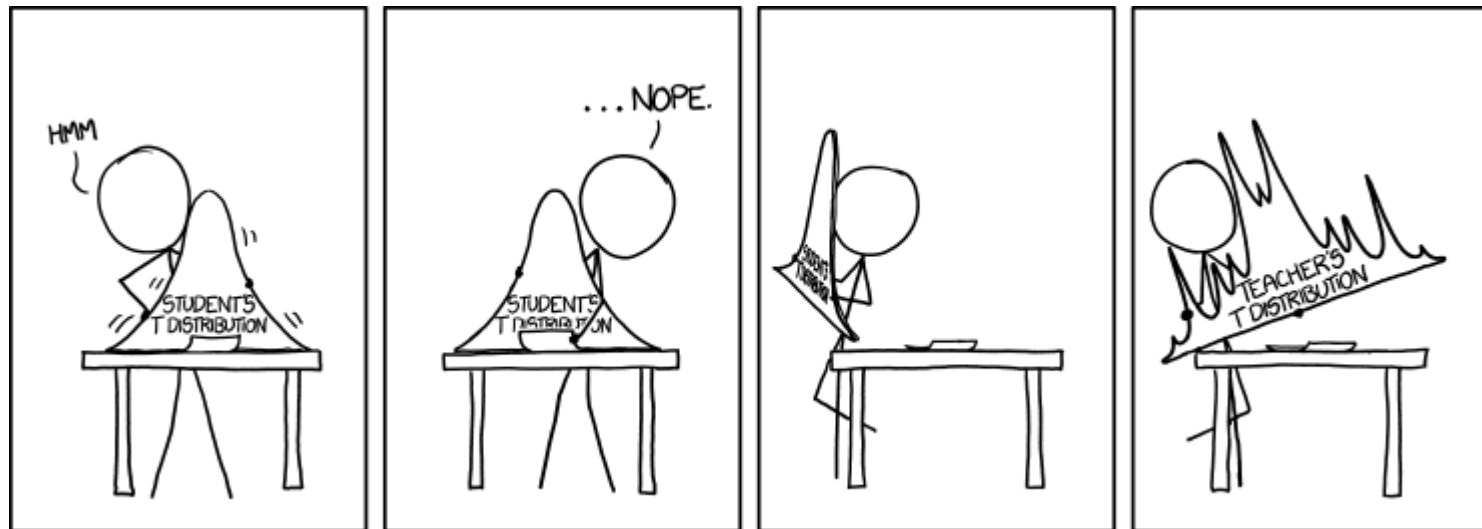
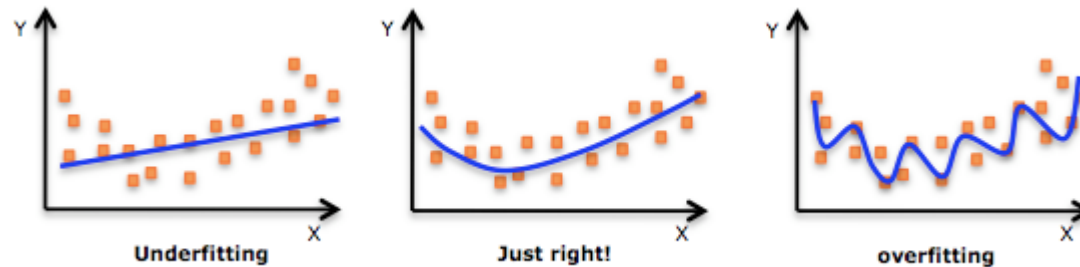
Increasing level of evidence for use of model in practice



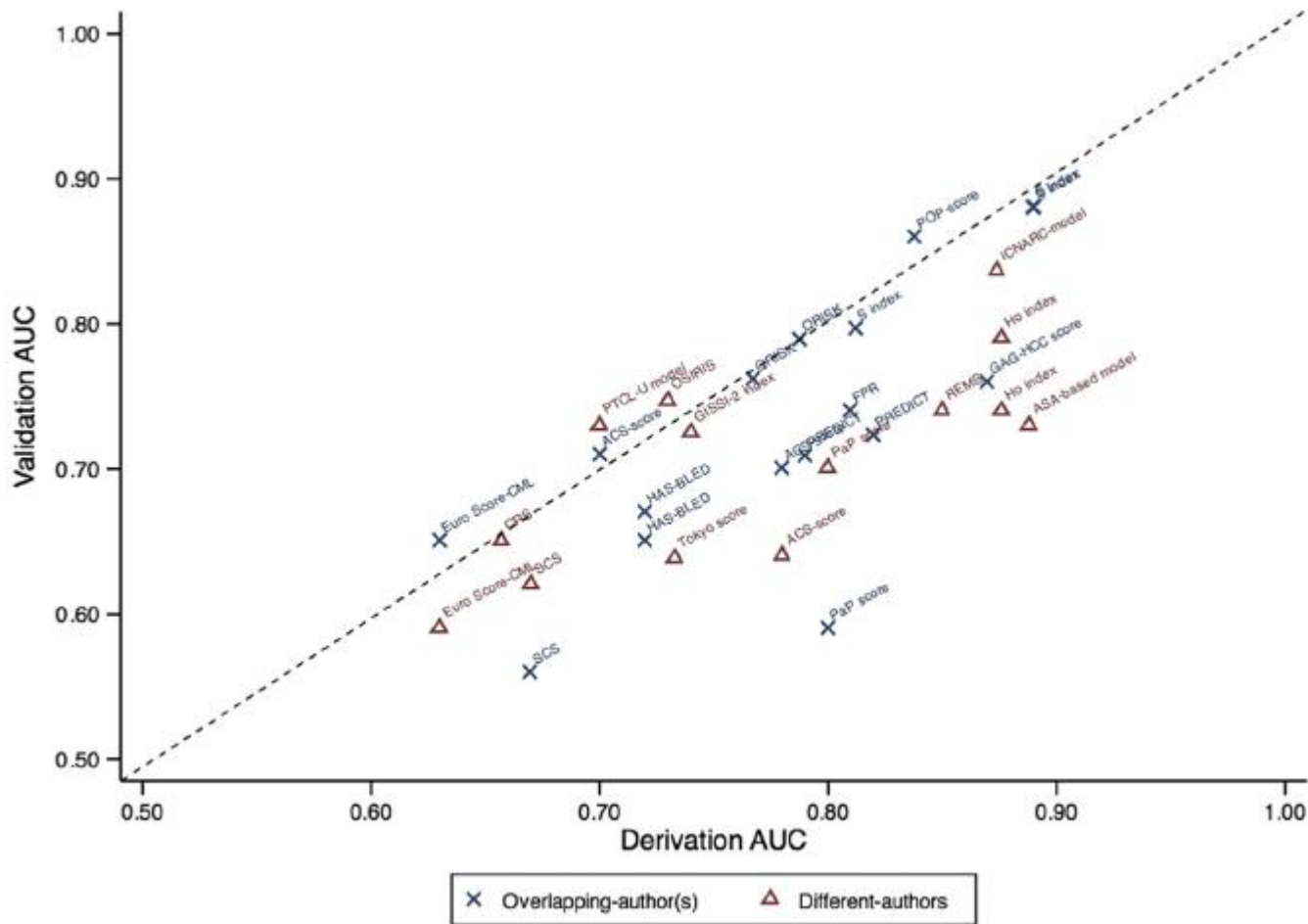
# Common pitfalls

# Lack of reproducibility

- Poor methodological & reporting standards
- Overfitting to data at hand



# Lack of transportability

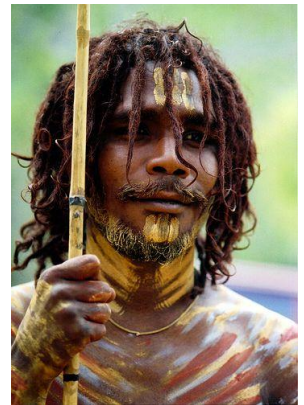
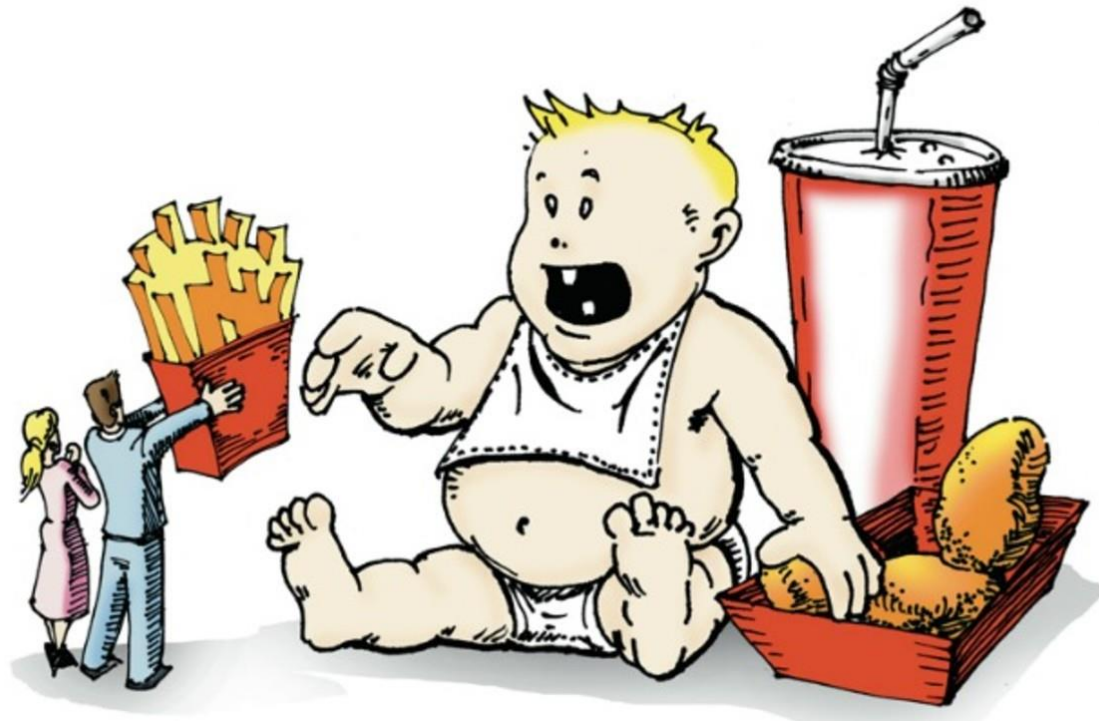


**Ref:** Siontis *et al.* External validation of new risk prediction models is infrequent and reveals worse prognostic discrimination. *Journal of Clinical Epidemiology*. 2014.



# Lack of transportability

- Missed important predictors
- Missed interaction terms & non-linear terms
- Poor measurement or modeling of relevant predictors



# Lack of transportability

- Differences in patient spectrum
- Differences in standards of care
- Differences in treatment standards



# Lack of (independent) validation



# Summarized

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

(and more often than not little attempt is made to address this issue)

- Poor quality of prognostic modelling studies
- Poor reproducibility
- Poor transportability
- Lack of external validation

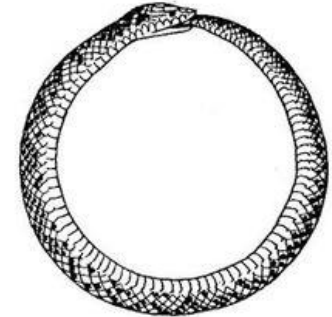
***“All models are wrong, but some are useful”***

George Box



# But wait... this is not the end

**There are numerous models for same target population and outcomes**



- > 150 models alike Framingham, SCORE, Qrisk
- > 100 models for brain trauma patients
- > 100 diabetes type 2 models
- > 60 models for breast cancer prognosis



# Numerous models for same target population + outcomes

*"Comparing risk prediction models should be routine when deriving a new model for the same purpose"* (Collins 2012)



*"Substantial work is needed to understand how competing prediction models compare and how they can best be applied to individualize care."* (Wessler 2015)



*"There is an excess of models predicting incident CVD in the general population. The usefulness of most of the models remains unclear."* (Damen 2016)

# Opportunities

**Evidence Synthesis**

Big Data

Machine Learning



UMC Utrecht

# Evidence synthesis

## Why?

- Improve estimation of prediction models
- Evaluate sources of variability in predictive performance
- Evaluate need for tailoring

## How?

- Synthesis of prognostic factors
- Synthesis of prediction models
- Synthesis of prediction model performance





# Evidence synthesis

Combining information on prognostic factors

**Concept:** Use previously published risk factor associations to update multivariable coefficients in “own” data set

Debray et al. *BMC Medical Research Methodology* 2012, **12**:121  
<http://www.biomedcentral.com/1471-2288/12/121>



TECHNICAL ADVANCE

Open Access

## Incorporating published univariable associations in diagnostic and prognostic modeling

Thomas P A Debray<sup>1\*</sup>, Hendrik Koffijberg<sup>1</sup>, Difei Lu<sup>2</sup>, Yvonne Vergouwe<sup>1,2</sup>,  
Ewout W Steyerberg<sup>2†</sup> and Karel G M Moons<sup>1†</sup>

STATISTICS IN MEDICINE  
*Statist. Med.* **19**, 141–160 (2000)

## PROGNOSTIC MODELS BASED ON LITERATURE AND INDIVIDUAL PATIENT DATA IN LOGISTIC REGRESSION ANALYSIS

E. W. STEYERBERG<sup>1\*</sup>, M. J. C. EUKEMANS<sup>1</sup>, J. C. VAN HOUWELINGEN<sup>2</sup>, K. L. LEE<sup>3</sup> AND  
J. D. F. HABBEMA<sup>1</sup>

<sup>1</sup>Center for Clinical Decision Sciences, Department of Public Health, Erasmus University, P.O. Box 1738, 3000 DR Rotterdam, The Netherlands

<sup>2</sup>Department of Medical Statistics, Leiden University, P.O. Box 9604, 2300 RC Leiden, The Netherlands

<sup>3</sup>Department of Community and Family Medicine, Duke University Medical Center, P.O. Box 3363, Durham, NC 27710, U.S.A.



# Evidence synthesis

Combining previously published prediction models

**Concept:** Use limited patient-level data at hand to combine and tailor previously published models

- Debray et al. *Statistics in Medicine* (2012) 31:23
- Debray et al. *Statistics in Medicine* (2014) 33:14
- Martin et al. *BMC Medical Research Methodology* (2017) 17:1



# Evidence synthesis

Combining previously published prediction models

### Diagnosis of Deep Vein Thrombosis

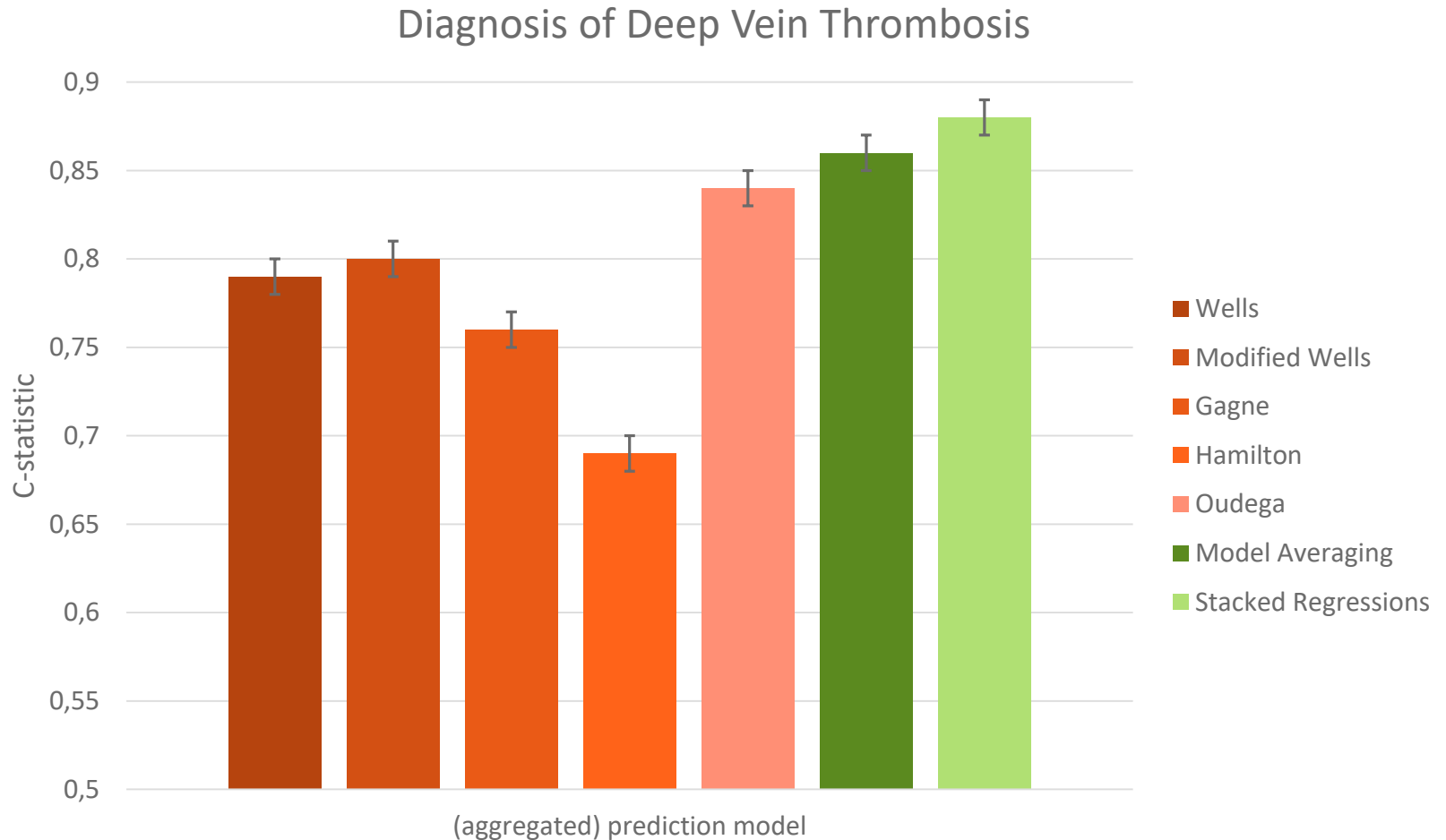


- Wells
- Modified Wells
- Gagne
- Hamilton
- Oudega
- Model Averaging
- Stacked Regressions



# Evidence synthesis

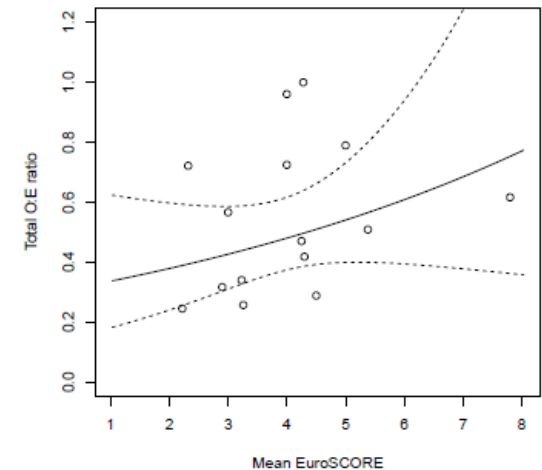
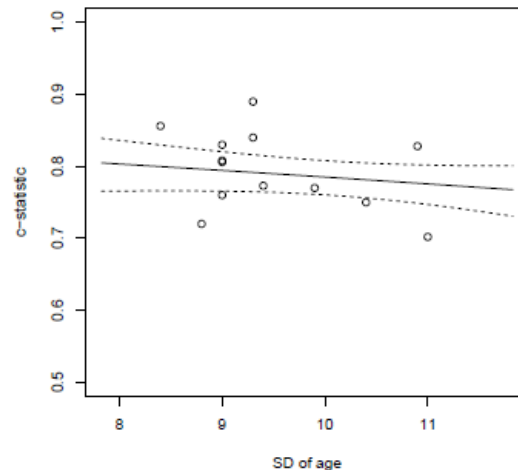
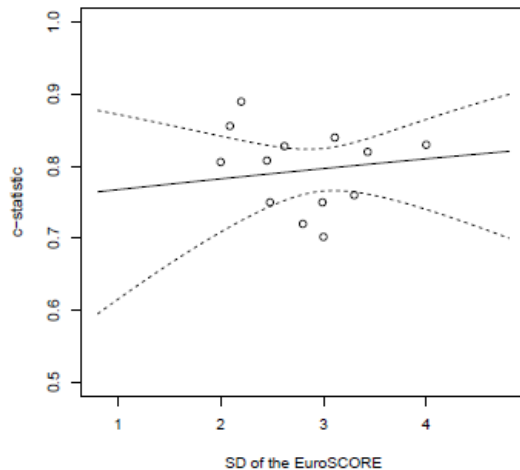
Combining previously published prediction models



# Evidence synthesis

Summarizing external validation study results

**Concept:** Systematically review external validation studies of a certain prediction model and summarize their results



**Ref:** Debray TPA, *et al.* A guide to systematic review and meta-analysis of prediction model performance. *BMJ* 2016 (Accepted for publication)



# Opportunities

Evidence Synthesis

**Big Data**

Machine Learning



UMC Utrecht

# The rise of big data

## What is 'big data'?

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

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

Example: QRISK2 was developed using e-health data from the QRESEARCH database using over 1.5 million patients (with over 95000 new cardiovascular events) from 355 randomly selected general practices



# Prediction research using big data

## Why do we need 'big data'?

- Development of better prediction models
  - Reduced risk of overfitting
  - Ability to address wider spectrum of patients
  - Ability to investigate more complex associations
- More extensive testing of model performance
  - Ability to externally validate across multiple settings (also upon model development)
  - Ability to investigate sources of poor or inconsistent model performance
  - Ability to assess usability of prediction models across different situations





# Prediction research using big data

## Prediction model development

Need to identify whether aggregation of IPD is justifiable, and how to adjust for heterogeneity.

- Allow for different baseline risks in each of the IPD studies or settings
- Investigate heterogeneity in predictor effects across studies or settings
- Implement a framework that uses internal-external cross-validation



# Internal-external cross-validation (IECV)

Statistics  
in Medicine



[Explore this journal >](#)

Research Article

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

Thomas P.A. Debray [✉](#), Karel G.M. Moons, Ikhlmaq Ahmed, Hendrik Koffijberg, Richard David Riley

First published: 11 January 2013 [Full publication history](#)

DOI: 10.1002/sim.5732 [View/save citation](#)

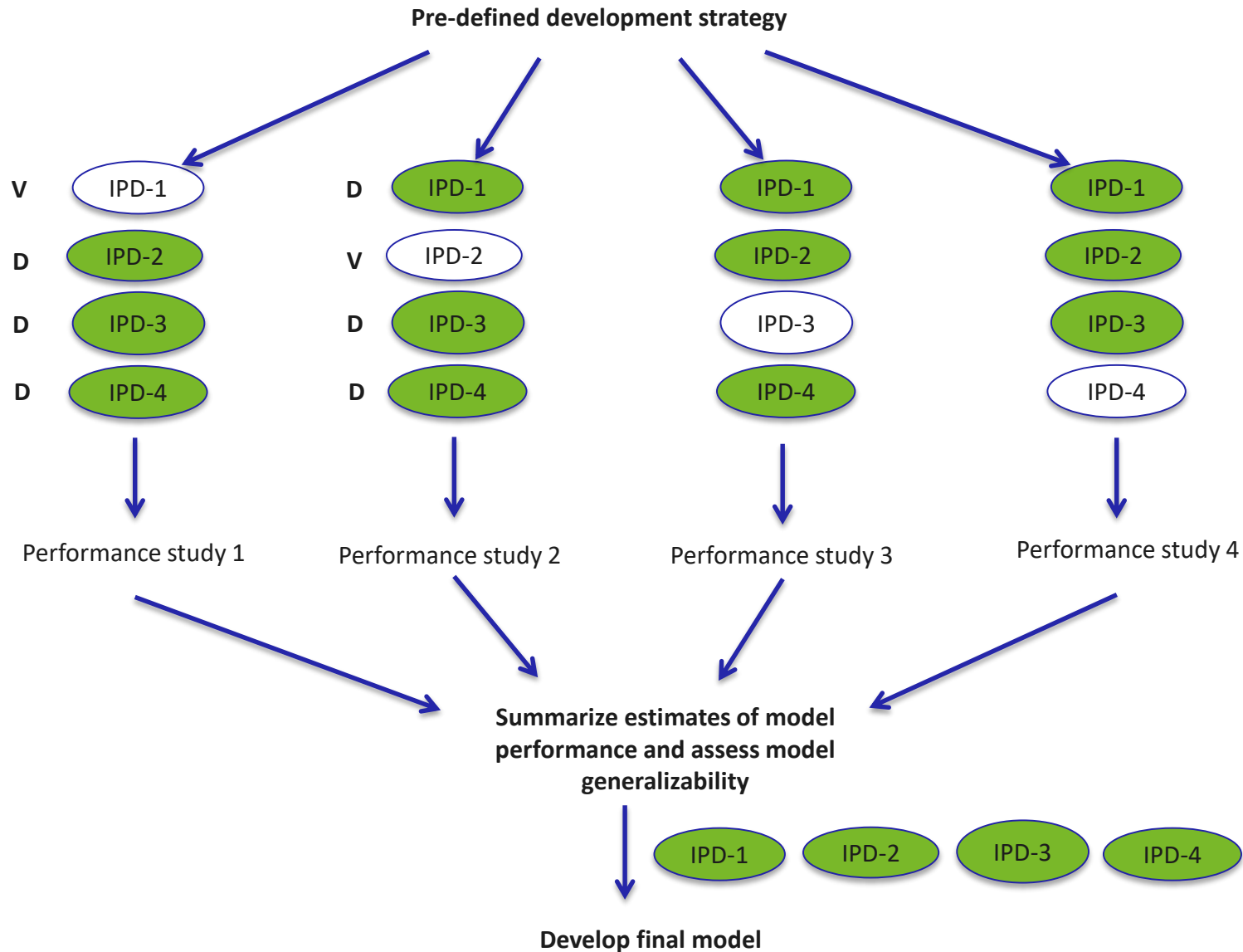
Cited by: 19 articles [Refresh](#) [Citing literature](#)



[View issue TOC](#)  
Volume 32, Issue 18  
15 August 2013  
Pages 3158-3180

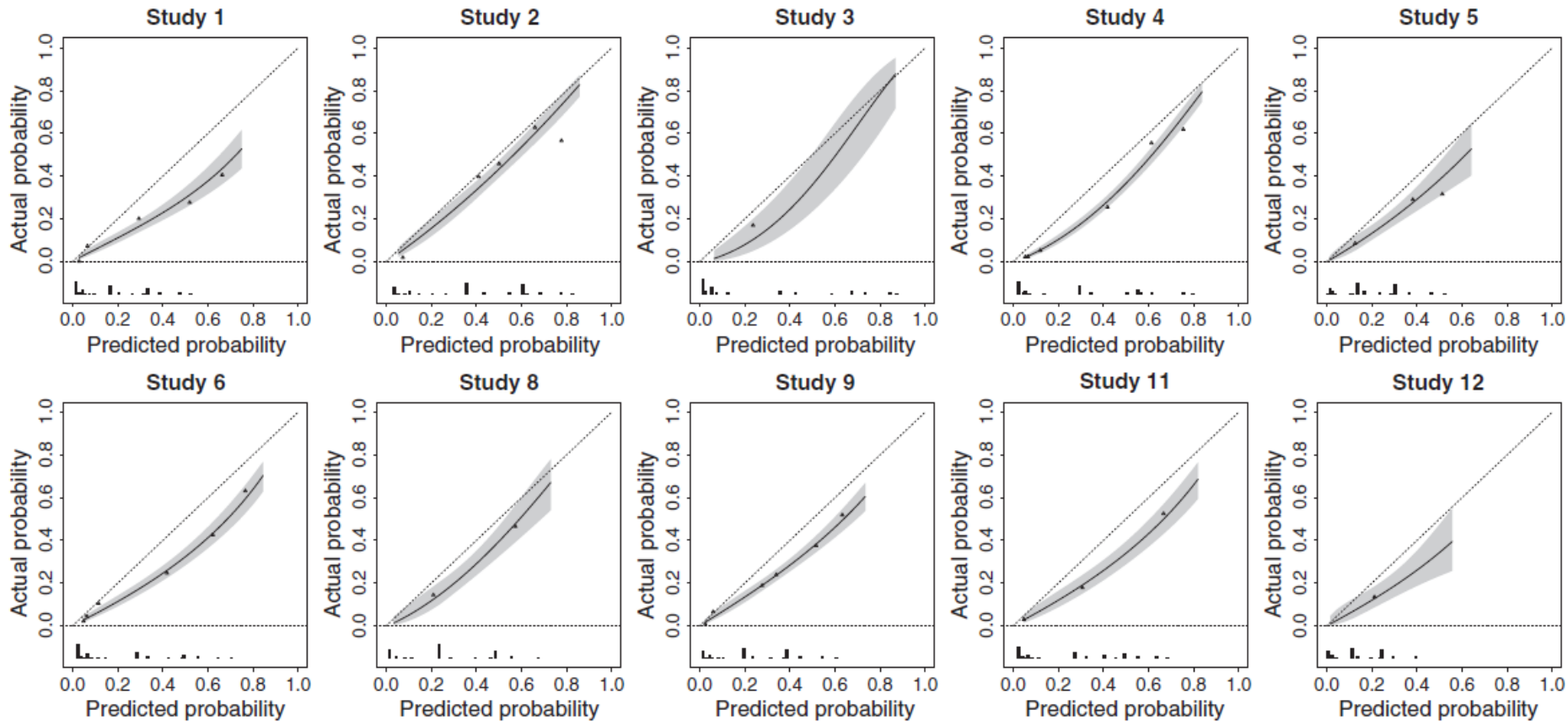
✉ Correspondence to: Thomas P. A. Debray, Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Stratenum 6.131, PO Box 85500, 3508GA Utrecht, The Netherlands.  
E-mail: T.Debray@umcutrecht.nl

# Internal-external cross-validation (IECV)



# Internal-external cross-validation (IECV)

The IECV approach allows for many external validations



# Assessing model performance

## Meta-analysis of performance estimates across different IPD sets

- A 'good' prediction model will have
  - satisfactory performance on average
  - little or no between-study heterogeneity in performance
- Need to summarize estimates of model performance...
  - To estimate likely performance in new studies
  - To calculate probability of "good" performance
  - To evaluate sources of between-study heterogeneity



# Meta-analysis of performance estimates



Journal of Clinical Epidemiology 69 (2016) 40–50

Journal of  
Clinical  
Epidemiology

Multivariate meta-analysis of individual participant data helped externally validate the performance and implementation of a prediction model

Kym I.E. Snell<sup>a</sup>, Harry Hua<sup>b</sup>, Thomas P.A. Debray<sup>c,d</sup>, Joie Ensor<sup>e</sup>,  
Maxime P. Look<sup>f</sup>, Karel G.M. Moons<sup>c,d</sup>, Richard D. Riley<sup>e,\*</sup>

<sup>a</sup>Public Health, Epidemiology and Biostatistics, School of Health and Population Sciences, Public Health Building, University of Birmingham, Edgbaston, Birmingham B15 2TT, UK

<sup>b</sup>School of Mathematics, Watson Building, University of Birmingham, Edgbaston, Birmingham B15 2TT, UK

<sup>c</sup>Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Str. 6.131, PO Box 85500, 3508 GA Utrecht, The Netherlands

<sup>d</sup>Dutch Cochrane Centre, University Medical Center Utrecht, Str. 6.131, PO Box 85500, 3508 GA Utrecht, The Netherlands

<sup>e</sup>Research Institute for Primary Care and Health Sciences, Keele University, Staffordshire ST5 5BG, UK

<sup>f</sup>Department of Medical Oncology, Erasmus MC Cancer Institute, Erasmus University Medical Center, PO Box 2040, 3000 CA Rotterdam, The Netherlands

Accepted 8 May 2015; Published online 16 May 2015

**Meta-analysis of prediction model performance across multiple studies: which scale helps ensure between-study normality for the C-statistic and calibration measures?**

Kym IE Snell<sup>1</sup>, Joie Ensor<sup>1</sup>, Thomas PA Debray<sup>2,3</sup>, Karel GM Moons<sup>2,3</sup>, Richard D Riley<sup>1</sup>

RESEARCH METHODS AND REPORTING

A guide to systematic review and meta-analysis of prediction model performance

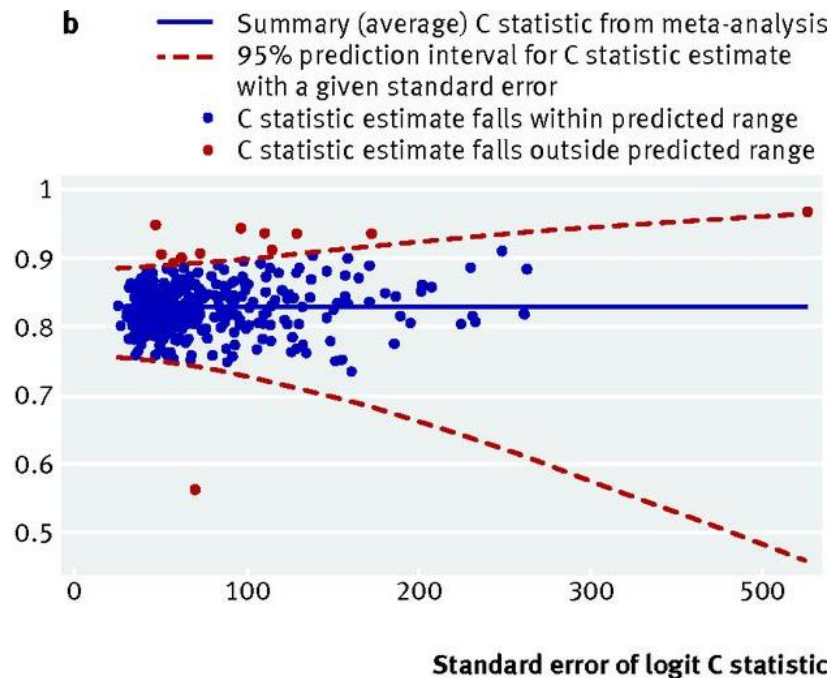
Thomas P A Debray,<sup>1,2</sup> Johanna A A G Damen,<sup>1,2</sup> Kym I E Snell,<sup>3</sup> Joie Ensor,<sup>3</sup> Lotty Hooft,<sup>1,2</sup> Johannes B Reitsma,<sup>1,2</sup> Richard D Riley,<sup>3</sup> Karel G M Moons<sup>1,2</sup>



click for updates

# Meta-analysis of performance estimates

## Evaluate model generalizability



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

Discrimination performance of QRISK2, across 364 general practice surgeries

**Ref:** Riley RD, *et al.* External validation of clinical prediction models using big datasets from e-health records or IPD meta-analysis: opportunities and challenges. *BMJ*. 2016;353:i3140.



# Meta-analysis of performance estimates

## Compare competing modeling strategies

- Choice of predictors
- Dealing with heterogeneity
- Non-linear effects
- Interaction terms

**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		
Calibration slope required	Minimum C statistic required	Strategy (1): Develop using logistic regression and implement with intercept estimated in external validation study	Strategy (2): Develop using logistic regression and implement with average study intercept taken from developed model	Strategy (3): Develop using logistic regression and implement with intercept taken from a study used in development data with a similar prevalence
0.9–1.1	0.70	0.027	0.037	0.037
0.8–1.2	0.70	0.146	0.158	0.156
0.9–1.1	0.65	0.427	0.413	0.409
0.8–1.2	0.65	0.728	0.712	0.707

*Abbreviation:* DVT, deep vein thrombosis.



# Meta-analysis of performance estimates

## Identify & address sources of heterogeneity

- Differences in patient spectrum
- Differences in baseline risk
- Differences in predictor effects

Facilitate tailoring of developed models!

### ORIGINAL ARTICLE

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

Thomas P.A. Debray<sup>a,\*</sup>, Yvonne Vergouwe<sup>b</sup>, Hendrik Koffijberg<sup>a</sup>, Daan Nieboer<sup>b</sup>,  
Ewout W. Steyerberg<sup>b,1</sup>, Karel G.M. Moons<sup>a,1</sup>

<sup>a</sup>*Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Str. 6.131, PO Box 85500, 3508GA Utrecht, The Netherlands*

<sup>b</sup>*Department of Public Health, Erasmus Medical Center, Rotterdam, The Netherlands*

Accepted 30 June 2014; Published online xxxx



GUIDELINES AND GUIDANCE

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

**Thomas P. A. Debray<sup>1,2\*</sup>, Richard D. Riley<sup>3</sup>, Maroeska M. Rovers<sup>4</sup>, Johannes B. Reitsma<sup>1,2</sup>, Karel G. M. Moons<sup>1,2</sup>, Cochrane IPD Meta-analysis Methods group<sup>¶</sup>**

**1** Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht, The Netherlands, **2** The Dutch Cochrane Centre, Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht, The Netherlands, **3** Research Institute for Primary Care and Health Sciences, Keele University, Staffordshire, The United Kingdom, **4** Radboud Institute for Health Sciences, Radboudumc Nijmegen, The Netherlands

<sup>¶</sup> Membership of the Cochrane IPD Meta-analysis Methods group is listed in the Acknowledgments.

\* [T.Debray@umcutrecht.nl](mailto:T.Debray@umcutrecht.nl)



CrossMark  
click for updates



# R package “metamisc”

**metamisc: Diagnostic and Prognostic Meta-Analysis**

Meta-analysis of diagnostic and prognostic modeling studies. Summarize estimates of diagnostic test accuracy and prediction model performance. Validate, update and combine published prediction models.

Version: 0.1.6  
Depends: R ( $\geq 2.10$ ), stats, graphics  
Imports: [metafor](#), [mvtnorm](#), [ellipse](#), [lme4](#)  
Suggests: [runjags](#), [rjags](#)  
Published: 2017-09-06  
Author: Thomas Debray [aut, cre], Valentijn de Jong [aut]  
Maintainer: Thomas Debray <thomas.debray at gmail.com>  
License: [GPL-2](#)  
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\\_0.1.6.tar.gz](#)  
Windows binaries: r-devel: [metamisc\\_0.1.6.zip](#), r-release: [metamisc\\_0.1.6.zip](#), r-oldrel: [metamisc\\_0.1.6.zip](#)  
OS X El Capitan binaries: r-release: [metamisc\\_0.1.6.tgz](#)  
OS X Mavericks binaries: r-oldrel: [metamisc\\_0.1.6.tgz](#)  
Old sources: [metamisc archive](#)

Linking:

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

# Opportunities

Evidence Synthesis

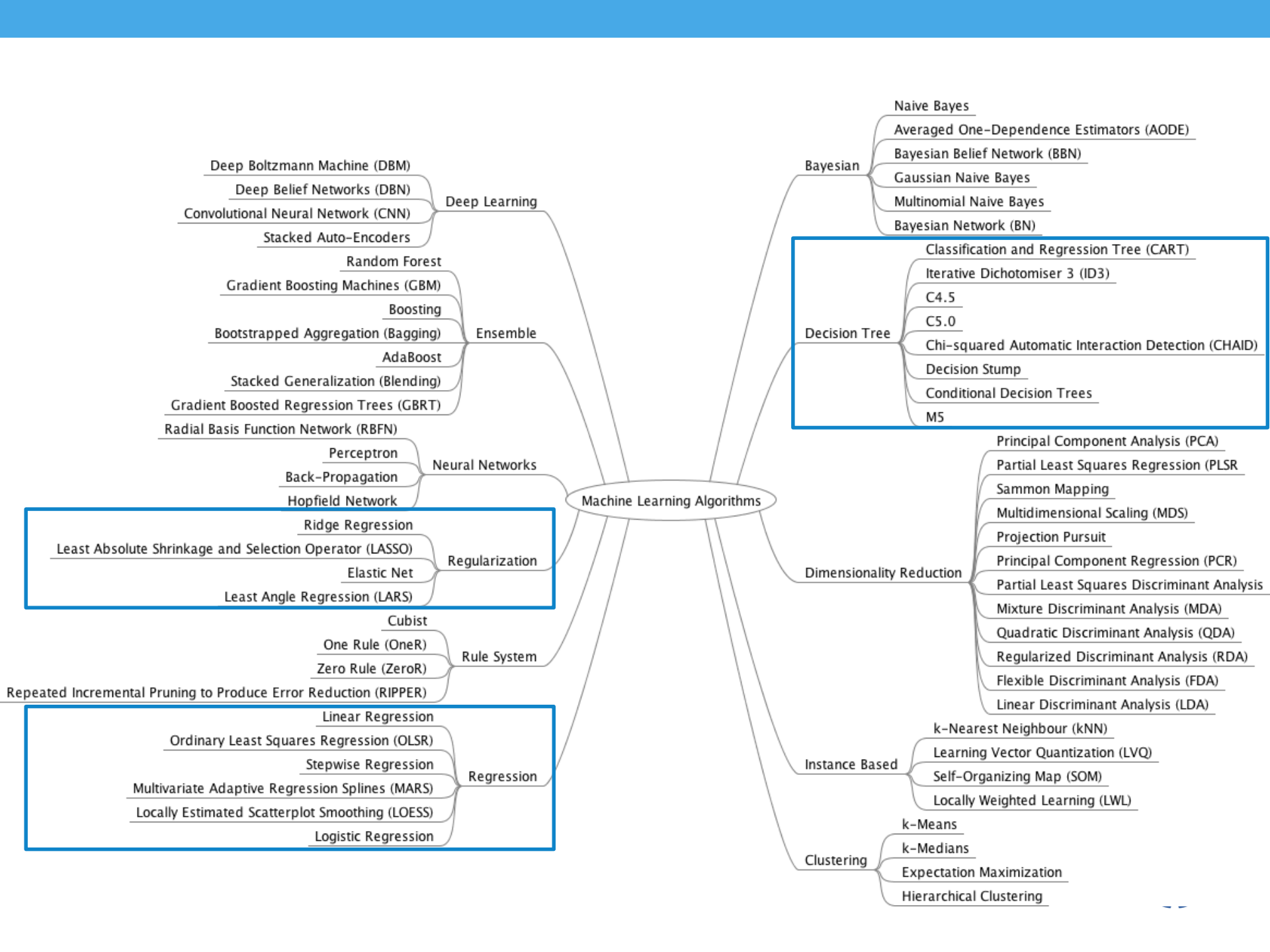
Big Data

**Machine Learning**



UMC Utrecht





Machine Learning Algorithms

- Deep Boltzmann Machine (DBM)
- Deep Belief Networks (DBN)
- Convolutional Neural Network (CNN)
- Stacked Auto-Encoders

Deep Learning

- Random Forest
- Gradient Boosting Machines (GBM)
- Boosting
- Bootstrapped Aggregation (Bagging)
- AdaBoost
- Stacked Generalization (Blending)
- Gradient Boosted Regression Trees (GBRT)

Ensemble

- Radial Basis Function Network (RBFN)
- Perceptron
- Back-Propagation
- Hopfield Network

Neural Networks

- Ridge Regression
- Least Absolute Shrinkage and Selection Operator (LASSO)
- Elastic Net
- Least Angle Regression (LARS)

Regularization

- Cubist
- One Rule (OneR)
- Zero Rule (ZeroR)
- Repeated Incremental Pruning to Produce Error Reduction (RIPPER)

Rule System

- Linear Regression
- Ordinary Least Squares Regression (OLSR)
- Stepwise Regression
- Multivariate Adaptive Regression Splines (MARS)
- Locally Estimated Scatterplot Smoothing (LOESS)
- Logistic Regression

Regression

- Naive Bayes
- Averaged One-Dependence Estimators (AODE)
- Bayesian Belief Network (BBN)
- Gaussian Naive Bayes
- Multinomial Naive Bayes
- Bayesian Network (BN)

Bayesian

- Classification and Regression Tree (CART)
- Iterative Dichotomiser 3 (ID3)
- C4.5
- C5.0
- Chi-squared Automatic Interaction Detection (CHAID)
- Decision Stump
- Conditional Decision Trees
- M5

Decision Tree

- Principal Component Analysis (PCA)
- Partial Least Squares Regression (PLSR)
- Sammon Mapping
- Multidimensional Scaling (MDS)
- Projection Pursuit
- Principal Component Regression (PCR)
- Partial Least Squares Discriminant Analysis
- Mixture Discriminant Analysis (MDA)
- Quadratic Discriminant Analysis (QDA)
- Regularized Discriminant Analysis (RDA)
- Flexible Discriminant Analysis (FDA)
- Linear Discriminant Analysis (LDA)

Dimensionality Reduction

- k-Nearest Neighbour (kNN)
- Learning Vector Quantization (LVQ)
- Self-Organizing Map (SOM)
- Locally Weighted Learning (LWL)

Instance Based

- k-Means
- k-Medians
- Expectation Maximization
- Hierarchical Clustering

Clustering

# Potential of Machine Learning

Machine Learning not widely implemented yet...

- Loss of transparency
- Performance gain often very limited



ELSEVIER

Journal of Clinical Epidemiology 65 (2012) 404–412

**Journal of  
Clinical  
Epidemiology**

Development and validation of clinical prediction models:  
Marginal differences between logistic regression, penalized maximum  
likelihood estimation, and genetic programming

Kristel J.M. Janssen<sup>a,\*</sup>, Ivar Siccama<sup>b</sup>, Yvonne Vergouwe<sup>a</sup>, Hendrik Koffijberg<sup>a</sup>, T.P.A. Debray<sup>a</sup>,  
Maarten Keijzer<sup>c</sup>, Diederick E. Grobbee<sup>a</sup>, Karel G.M. Moons<sup>a</sup>

<sup>a</sup>Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, P.O. Box 85500, 3508 AB Utrecht, The Netherlands

<sup>b</sup>Department of Neurology, Erasmus Medical Center, Rotterdam, The Netherlands

<sup>c</sup>Pegasystems Benelux, Amsterdam, The Netherlands

Accepted 9 August 2011; Published online 02 January 2012



# Potential of Machine Learning

With the rise of big data, the appeal of machine learning is increasing.

Key strengths

- Handling enormous numbers of predictors
- Modeling highly interactive and nonlinear effects





# Potential of Machine Learning

Promising areas of application

- Analysis of unstructured data
  - Text (e.g. medical records)
  - Images (e.g. CT, MRI, ...)
- Analysis of high velocity data
  - Brain signals (e.g. restoration of motor control)
  - Wearable devices
  - Social media
- Diagnosis
  - Generation of differential diagnoses
  - Suggestion of high-value tests



# Reasons to be optimistic?

