# Clinical prediction models in the age of artificial intelligence and big data

Ewout Steyerberg

*Professor of Clinical Biostatistics and Medical Decision Making*

<E.Steyerberg@ErasmusMC.nl /
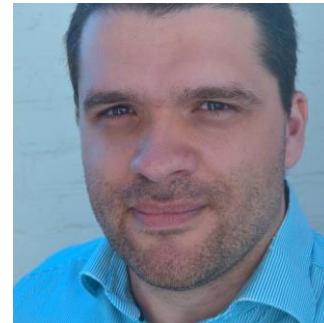
E.W.Steyerberg@LUMC.nl >

Basel, Nov 1 2019

LU
MC

LEIDEN UNIVERSITY MEDICAL CENTER

Erasmus MC
University Medical Center Rotterdam

# Thanks to co-workers; no COI

- LUMC: Maarten van Smeden

- Leuven: Ben van Calster

Both provided many of the slides shown

# Main question

Where does Big Data / machine learning (ML) / artificial intelligence (AI) assist us in prediction research?

- Strengths and weaknesses of Big Data initiatives

- Consider links between classical statistical approaches, ML, AI for prediction

# Prediction models; what for?

- Understanding nature:
  relative risks of different predictors

- Predicting outcomes:
  absolute risk by combinations of predictors

## To Explain or to Predict?

**Galit Shmueli**

# Traditional regression modeling

Can well be used for explanation and prediction

# Prediction models

- Diagnosis
  - Imaging findings, e.g. abnormal CT scan in trauma
  - Clinical condition, e.g. serious infection
  - …
- Prognosis
  - Mortality, e.g. < 30 days, over time, …
  - …

# Prognostic / predictive models

Prognostic modeling

y ~ X          Prognostic factors

y ~ Tx         Treatment effect

y ~ X + Tx     Covariate adjusted tx effect
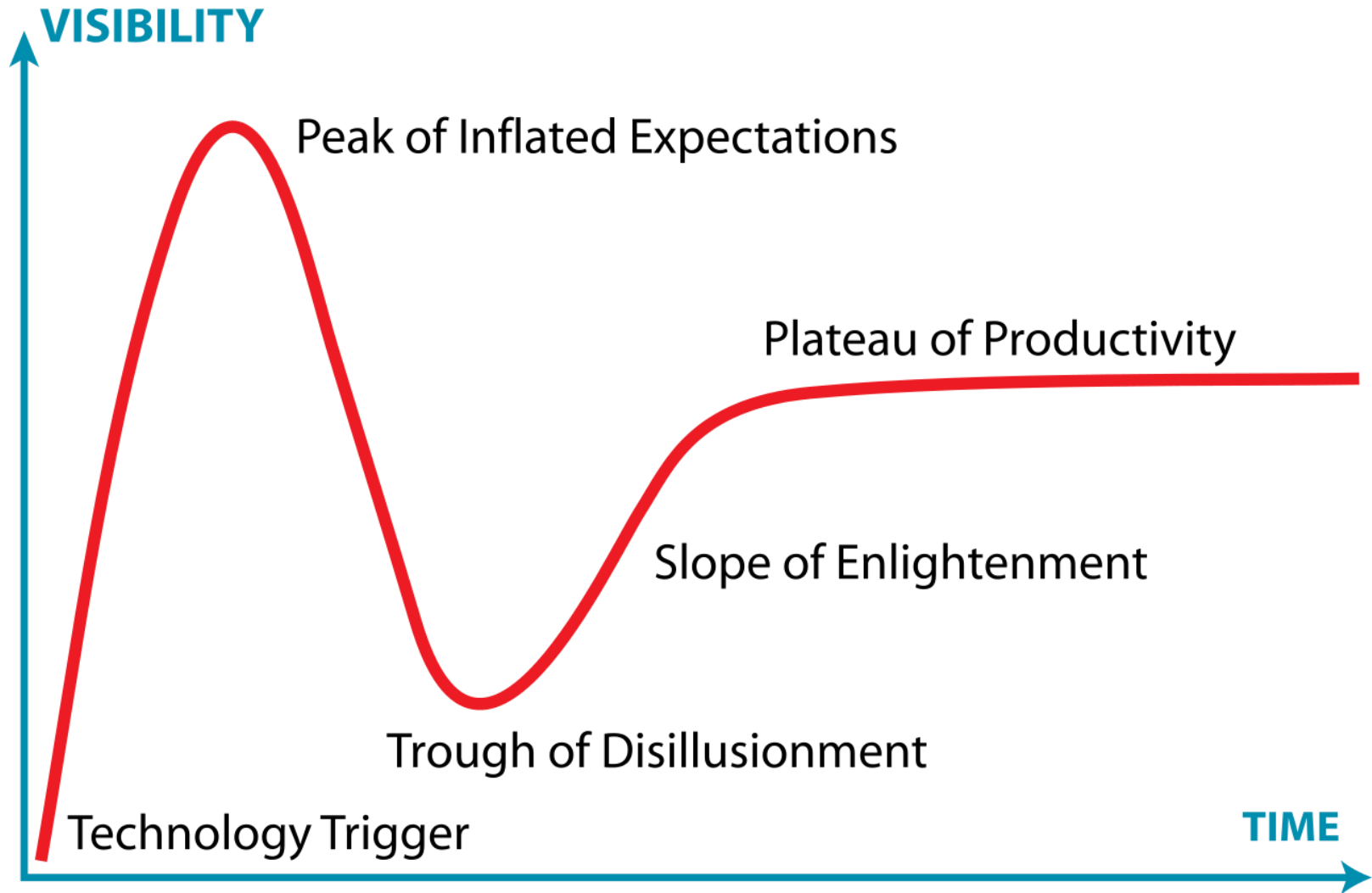
Predictive modeling

y ~ X * Tx     Predictive factors for differential tx effect

# Opportunities in medical prediction

- More data
  - larger N
  - more variables
- More detail
  - biomarkers / omics / imaging / eHealth
- Novel methods
  - ML / AI / ..
  - Statistical methods
    - Dynamic prediction
    - Testing procedures for high dimensional data
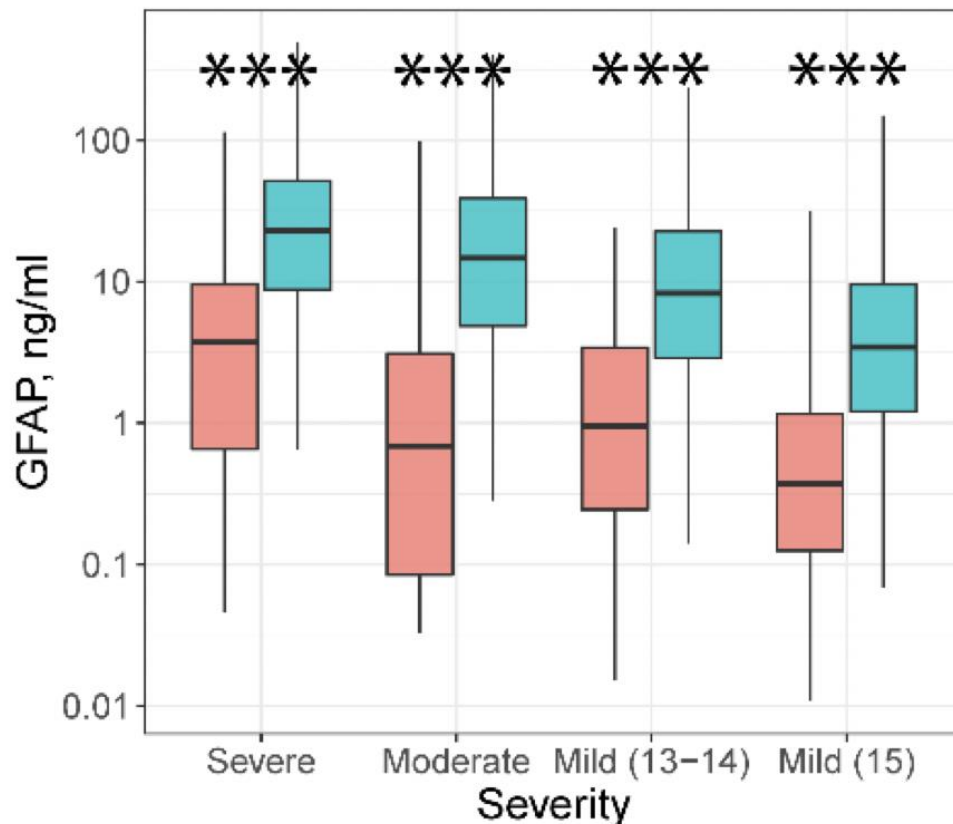    - …

# Hype

# Examples

- Biomarkers

- Imaging

- Omics

# Positive example 1

- Biomarkers in diagnosing head trauma
  - Mild: AUC 0.89 [0.87-0.90] vs clinical 0.84 [0.83-0.86]

# Positive example 2

- MRI Imaging in diagnosing prostate cancer

**External Validation and Comparison of Prostate Cancer Risk Calculators Incorporating Multiparametric Magnetic Resonance Imaging for Prediction of Clinically Significant Prostate Cancer.**

Saba K[1], Wettstein MS[1,2], Lieger L[1], Hötker AM[3], Donati OF[3], Moch H[4], Ankerst DP[5], Poyet C[1], Sulser T[1], Eberli D[1], Mortezavi A[1].

- MRI-PCa-RCs AUC **0.83 to 0.85** vs
  PCa-RCs AUC **0.69 to 0.74**

# Positive example 3



Omics revolution

Metabolomics — Metabolites (A, B, C, D)
Fluxomics — Fluxes ($V_{AB}$, $V_{CD}$)
Proteomics — Proteins
Transcriptomics — mRNA
Genomics — DNA

System biology

Integrative physiology

System medicine

System pharmacology

Regenerative medicine

Integrated biomarkers

Human disease
  Prediction
  Diagnostics
  Treatment efficacy

# Positive example 3

- Omics in diagnosing … / predicting … ??

- Because omics →
          clinical characteristics →
                              outcome?

# Examples

- Biomarkers
- Imaging
- Omics

- ML / AI

# Success of ML / AI

# Non-exhaustive list

Gaming

Natural Language Processing (Siri etc)

Fraud detection

Shoplifting

Object recognition (e.g. for driverless cars)

Facial recognition

Traffic predictions (e.g. Waze app)

Electrical load forecasting

(Social) media and advertising (people you may know, movie suggestions, )
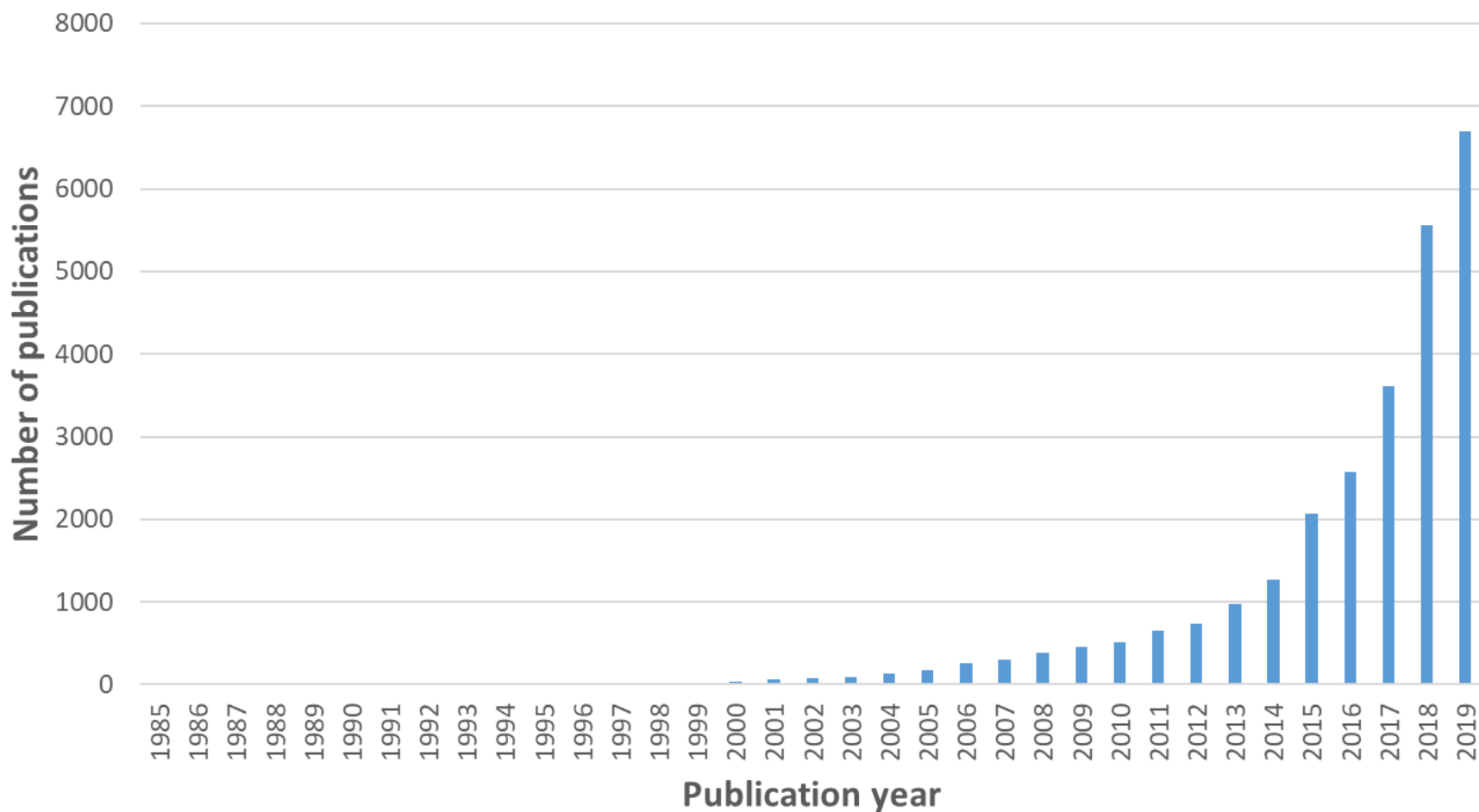
Spam filtering

Search engines (e.g. Google PageRank)

Handwriting recognition

# Popularity skyrocketing



"machine learning" in Medline database

# IBM Watson winning Jeopardy! (2011)

# IBM Watson for oncology

STAT+

**IBM's Watson supercomputer recommended 'unsafe and incorrect' cancer treatments, internal documents show**

By CASEY ROSS @caseymross and IKE SWETLITZ / JULY 25, 2018

Internal IBM documents show that its Watson supercomputer often spit out erroneous cancer treatment advice and that company medical specialists and customers identified "multiple examples of unsafe and incorrect treatment recommendations" as IBM was promoting the product to hospitals and physicians around the world.

https://bit.ly/2LxiWGj

# Evidence

- Cochrane: "We searched for RCTs and found 20 among … papers"

- Dr Watson: "We searched 4 Million webpages in 1 second"

# Five myths

1. Big Data will resolve the problems of small data
2. ML/AI is very different from classical modeling
3. Deep learning is relevant for all medical prediction problems
4. ML / AI is better than classical modeling for medical prediction problems
5. ML / AI leads to better generalizability

# Myth 1: Big Data will resolve the problems of small data

# High-performance medicine: the convergence of human and artificial intelligence
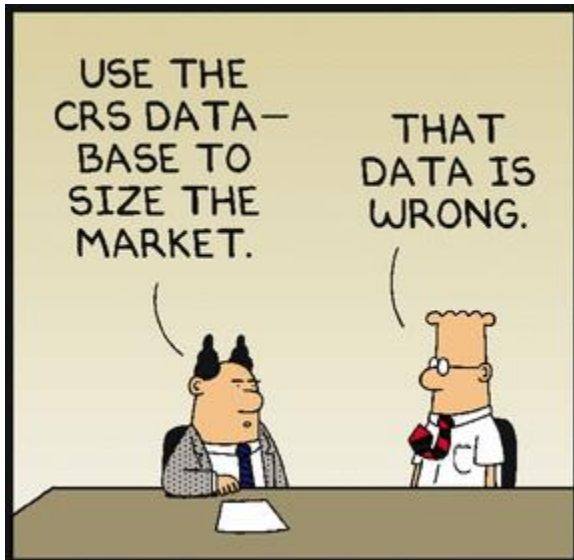
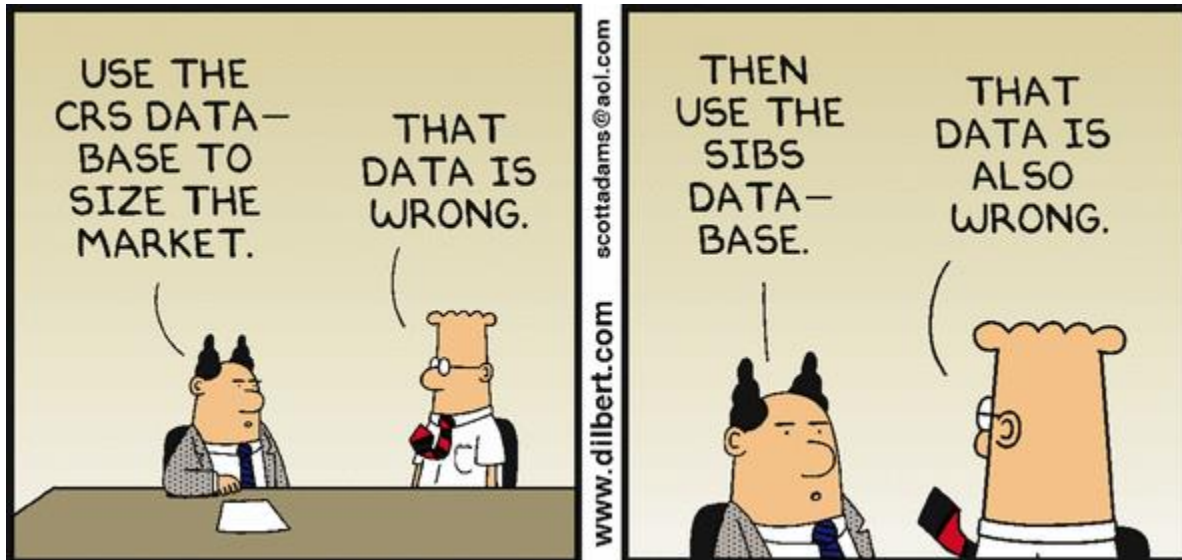Eric J. Topol ✉

## Abstract

The use of artificial intelligence, and deep-learning in particular, has been enabled by the use of big data, along with markedly enhanced computing power and cloud storage, across all sectors.

In medicine, this is beginning to have an impact ...
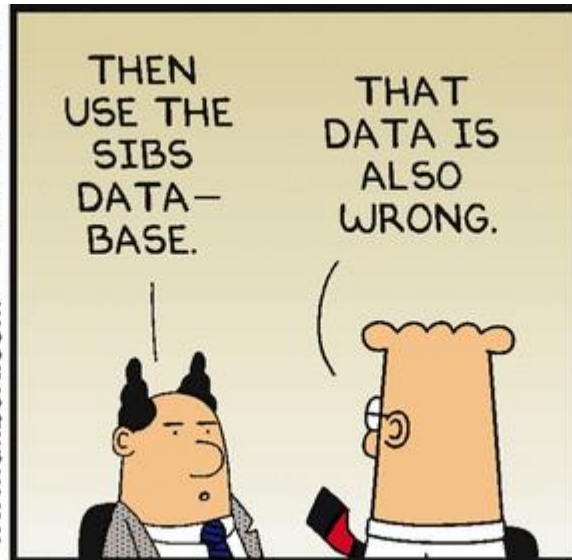
Do you have a clear research question?
Do you have data that help you answer the question?
What is the quality of the data?

Do you have a clear research question?
Do you have data that help you answer the question?
What is the quality of the data?

Do you have a clear research question?
Do you have data that help you answer the question?
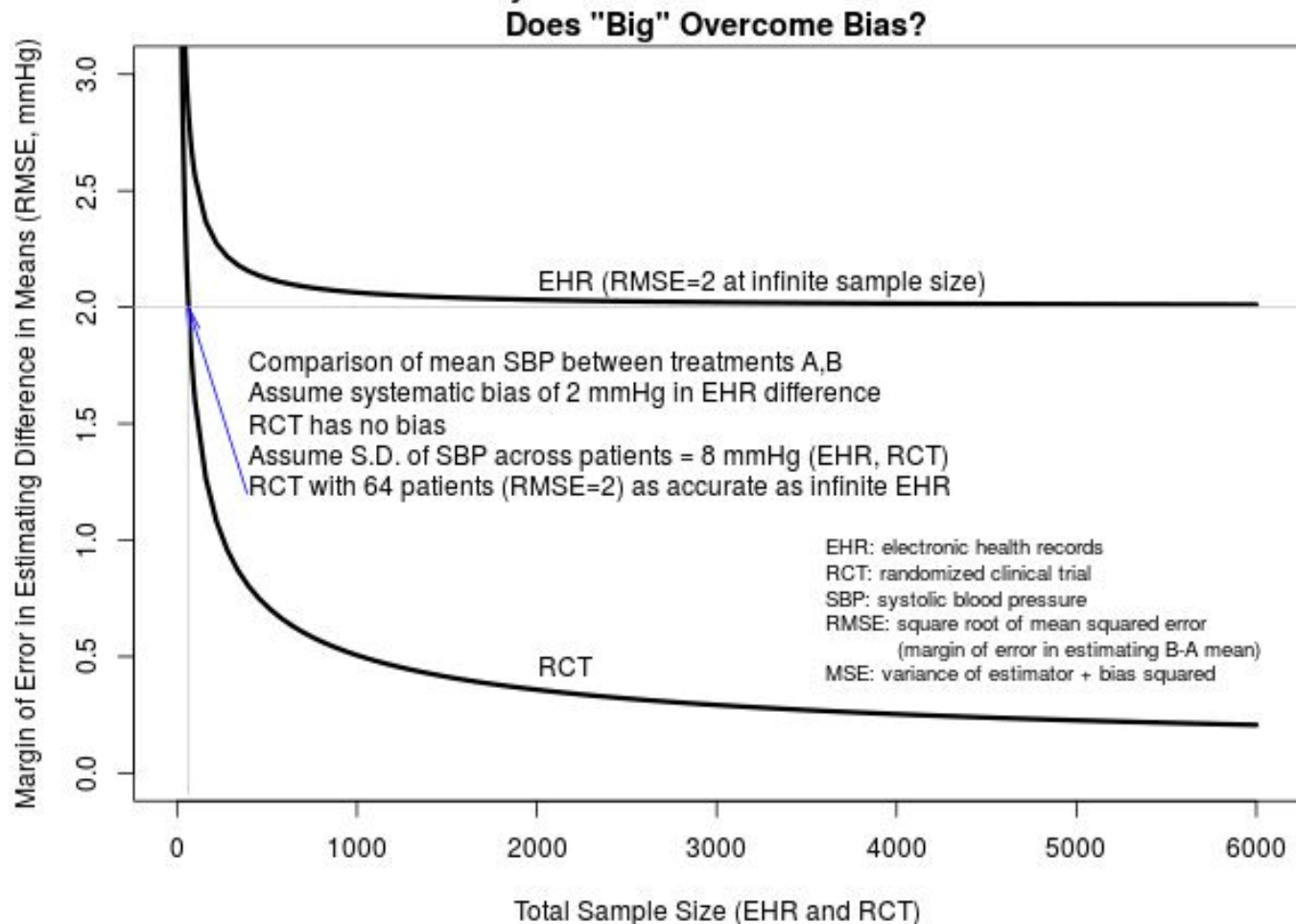What is the quality of the data?

# Big Data, Big Errors



Frank Harrell @f2harrell · 23 jun. 2017
Example: RCT randomizing 64 patients as accurate as infinitely large EHR: fharrell.com/2017/06/ehrs-a... #StatThink #RCT #EHR #BigData

Does "Big" Overcome Bias?

EHR (RMSE=2 at infinite sample size)

Comparison of mean SBP between treatments A,B
Assume systematic bias of 2 mmHg in EHR difference
RCT has no bias
Assume S.D. of SBP across patients = 8 mmHg (EHR, RCT)
RCT with 64 patients (RMSE=2) as accurate as infinite EHR

EHR: electronic health records
RCT: randomized clinical trial
SBP: systolic blood pressure
RMSE: square root of mean squared error
(margin of error in estimating B-A mean)
MSE: variance of estimator + bias squared

RCT

Margin of Error in Estimating Difference in Means (RMSE, mmHg)

Total Sample Size (EHR and RCT)

# Myth 2: ML/AI is very different from classical modeling

# "Everything is ML"



Scott H. Hawley
@drscotthawley

Replying to @JuliaHCox, @mikarv and @GSCollins

Logistic regression IS machine learning.

4:17 pm · 17 Feb 2019 · Twitter for iPhone

https://bit.ly/2IEVn33

# Two cultures

## Statistical Modeling: The Two Cultures

**Leo Breiman**

# Traditional Statistics vs Machine Learning

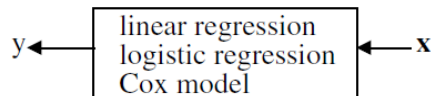## Statistical Modeling: The Two Cultures

**Leo Breiman**

**The Data Modeling Culture**

The analysis in this culture starts with assuming a stochastic data model for the inside of the black box. For example, a common data model is that data are generated by independent draws from
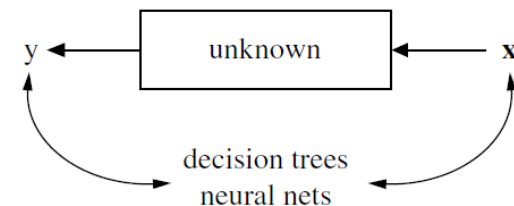
response variables $= f$(predictor variables, random noise, parameters)

The values of the parameters are estimated from the data and the model then used for information and/or prediction. Thus the black box is filled in like this:
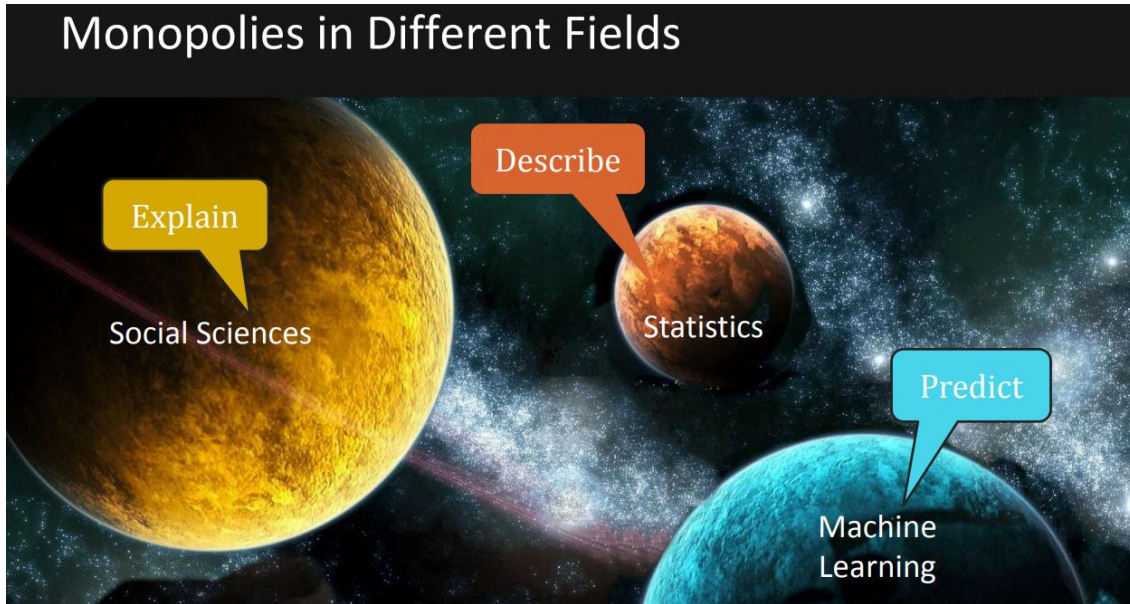
y ← [ linear regression / logistic regression / Cox model ] ← x

**The Algorithmic Modeling Culture**

The analysis in this culture considers the inside of the box complex and unknown. Their approach is to find a function $f(\mathbf{x})$—an algorithm that operates on $\mathbf{x}$ to predict the responses $\mathbf{y}$. Their black box looks like this:

y ← [ unknown ] ← x

decision trees
neural nets

# Traditional Statistics vs Machine Learning

# Example of exaggerating contrasts

RESEARCH ARTICLE

Machine learning models in electronic health records can outperform conventional survival models for predicting patient mortality in coronary artery disease

Andrew J. Steele[1]*, Spiros C. Denaxas[2], Anoop D. Shah[2,3], Harry Hemingway[2], Nicholas M. Luscombe[1,4,5]

**Table 1. The 27 expert-selected predictors used.**

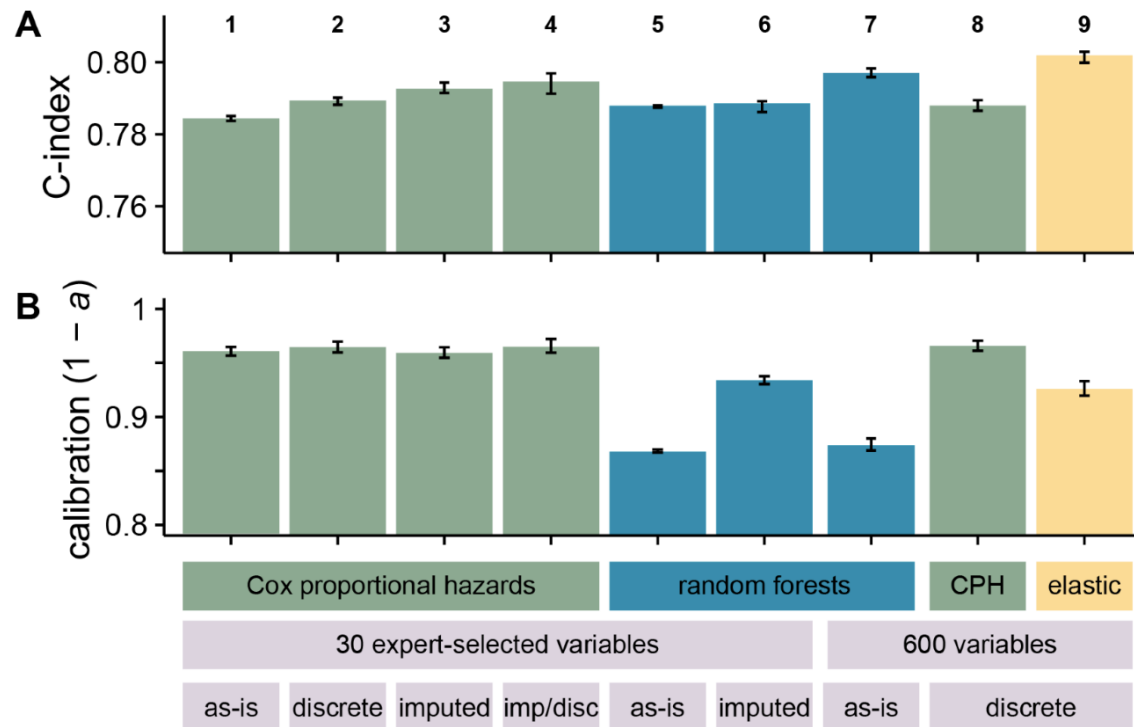| Category | Prognostic factors |
| --- | --- |
| Sociodemographic characteristics | Age, gender, most deprived quintile |
| CVD diagnosis and severity | SCAD subtype (stable angina, unstable angina, STEMI, NSTEMI, other CHD), PCI in last six months, CABG in last six months, previous/recurrent MI, use of nitrates |
| CVD risk factors | Smoking status (current, ex, never), hypertension, diabetes mellitus, total cholesterol, HDL |
| CVD comorbidities | Heart failure, peripheral arterial disease, atrial fibrillation, stroke |
| Non-CVD comorbidities | Chronic kidney disease, chronic obstructive pulmonary disease, cancer, chronic liver disease |
| Psychosocial characteristics | Depression at diagnosis, anxiety at diagnosis |
| Biomarkers | Heart rate, creatinine, white cell count, haemoglobin |

# Predicting mortality – the results



Fig 1. **Overall discrimination and calibration performance for the different models and datasets used.** (A) shows discrimination (C-

Elastic net, 586 ('600') variables: $c$=0.801

Traditional Cox, 27 ('30') expert-selected variables: $c$=0.793

# Predicting mortality – the media



AI beats doctors at predicting heart disease deaths

4 SEPTEMBER 2018   HUMAN BIOLOGY   HEALTH AND AGEING   NEWS

AI NEWS RESEARCH —

Artificial Intelligence beats doctors at predicting heart disease deaths

BY SHACK15 – 5 SEPTEMBER, 2018

**Science**Daily®

Your source for the latest research news

SD   Health ▾   Tech ▾   Enviro ▾   Society ▾   Quirky ▾

Science News                                   from research organization

AI beats doctors at predicting heart disease deaths
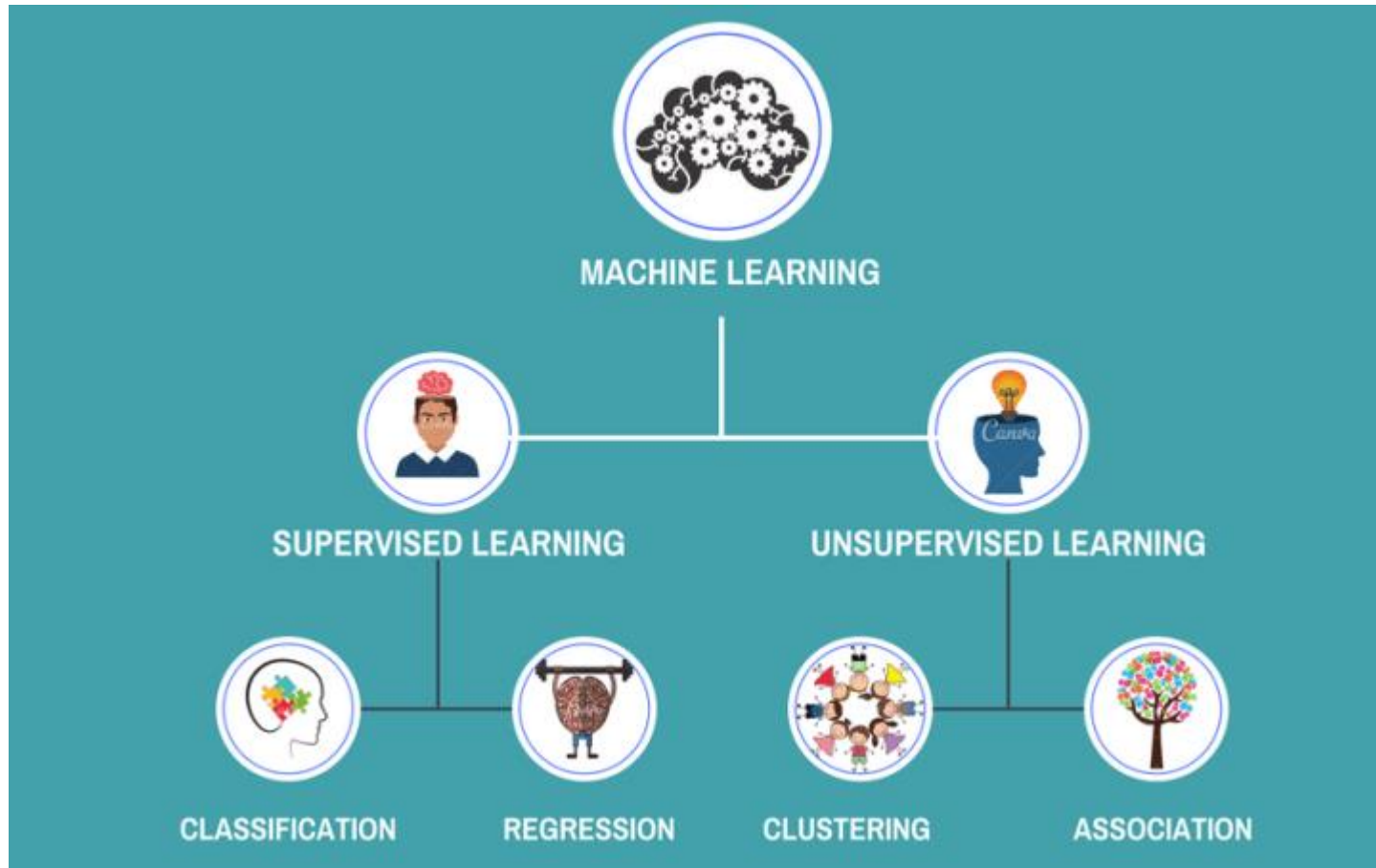
Date:   September 4, 2018

# ML refers to a culture, not to methods

- Substantial **overlap methods** used by both cultures
- Substantial **overlap analysis goals**
- Attempts to separate the two frequently result in **disagreement**

**Pragmatic approach:**

"ML" refers to models roughly outside of the traditional regression types of analysis:
trees, SVMs, neural networks, boosting etc.

# Machine learning: simple overview

# Myth 3: Deep learning is relevant for all medical prediction
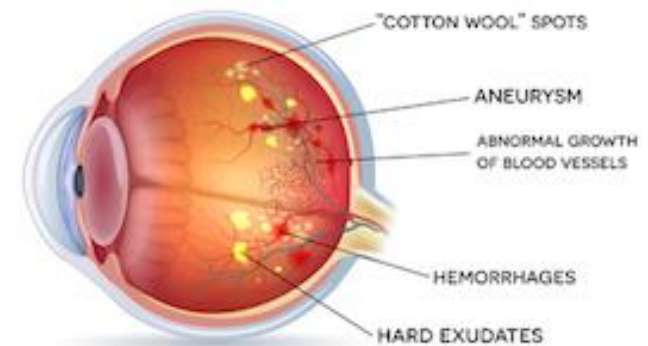
# Example: retinal disease

## Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs

Varun Gulshan, PhD; Lily Peng, MD, PhD; Marc Coram, PhD; Martin C. Stumpe, PhD; Derek Wu, BS; Arunachalam Narayanaswamy, PhD; Subhashini Venugopalan, MS; Kasumi Widner, MS; Tom Madams, MEng; Jorge Cuadros, OD, PhD; Ramasamy Kim, OD, DNB; Rajiv Raman, MS, DNB; Philip C. Nelson, BS; Jessica L. Mega, MD, MPH; Dale R. Webster, PhD

Diabetic retinopathy

Deep learning (= Neural network)

- 128,000 images

- Transfer learning (preinitialization)

- Sensitivity and specificity > .90
    - Estimated from training data

"COTTON WOOL" SPOTS

ANEURYSM

ABNORMAL GROWTH OF BLOOD VESSELS

HEMORRHAGES

HARD EXUDATES

Gulshan et al, JAMA, 2016, 10.1001/jama.2016.17216;
Picture retinopathy: https://bit.ly/2kB3X2w AS

# Example: lymph node metastases

JAMA | Original Investigation

## Diagnostic Assessment of Deep Learning Algorithms for Detection of Lymph Node Metastases in Women With Breast Cancer

Babak Ehteshami Bejnordi, MS; Mitko Veta, PhD; Paul Johannes van Diest, MD, PhD; Bram van Ginneken, PhD;
Nico Karssemeijer, PhD; Geert Litjens, PhD; Jeroen A. W. M. van der Laak, PhD; and the CAMELYON16 Consortium

Deep learning competition

But:
- 390 teams signed up, 23 submitted
- "Only" 270 images for training
- Test AUC range: 0.56 to 0.99

| Codename[b] | Task 1: Metastasis Identification FROC Score (95% CI)[c] | Task 2: Metastases Classification AUC (95% CI)[c] |
|---|---|---|
| HMS and MIT II | 0.807 (0.732-0.889) | 0.994 (0.983-0.999) |
| HMS and MGH III | 0.760 (0.692-0.857) | 0.976 (0.941-0.999 |
| HMS and MGH I | 0.596 (0.578-0.734) | 0.964 (0.928-0.989) |
| VISILAB II | 0.116 (0.063-0.177) | 0.651 (0.549-0.742) |
| Anonymous I | 0.097 (0.049-0.158) | 0.628 (0.530-0.717) |
| Laboratoire d'Imagerie Biomédicale I | 0.120 (0.079-0.182) | 0.556 (0.434-0.654) |

3. Deep learning is relevant for all medical prediction problems

**NO: Deep learning excels in visual tasks**

# Myth 4:   ML / AI is better than classical modeling for medical prediction

# Reviewer #2,
# van Smeden submission 2019

used in this paper. Second, since the prediction performance of logistic regression models is often inferior to those of powerful machine learning algorithms such as random forest or boosting, focussing logistic regression models only can be boring. The detailed comments are given below.

# REVIEW

# A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models

Evangelia Christodoulou[a], Jie Ma[b], Gary S. Collins[b,c], Ewout W. Steyerberg[d], Jan Y. Verbakel[a,e,f], Ben Van Calster[a,d,*]

# Poor methods and unclear reporting

What was done about missing data? 45% fully unclear, 100% poor or unclear

How were continuous predictors modeled? 20% unclear, 25% categorized

How were hyperparameters tuned? 66% unclear, 19% tuned with information

How was performance validated? 68% unclear or biased approach

Was accuracy of risk estimates checked? 79% not at all


Further observations:
- Prognosis: time horizon often ignored
- Patients matched on variables used a predictors
- 99% of patients excluded from modeling to obtain a balanced dataset
- First and last percentile of continuous predictors replaced with mean
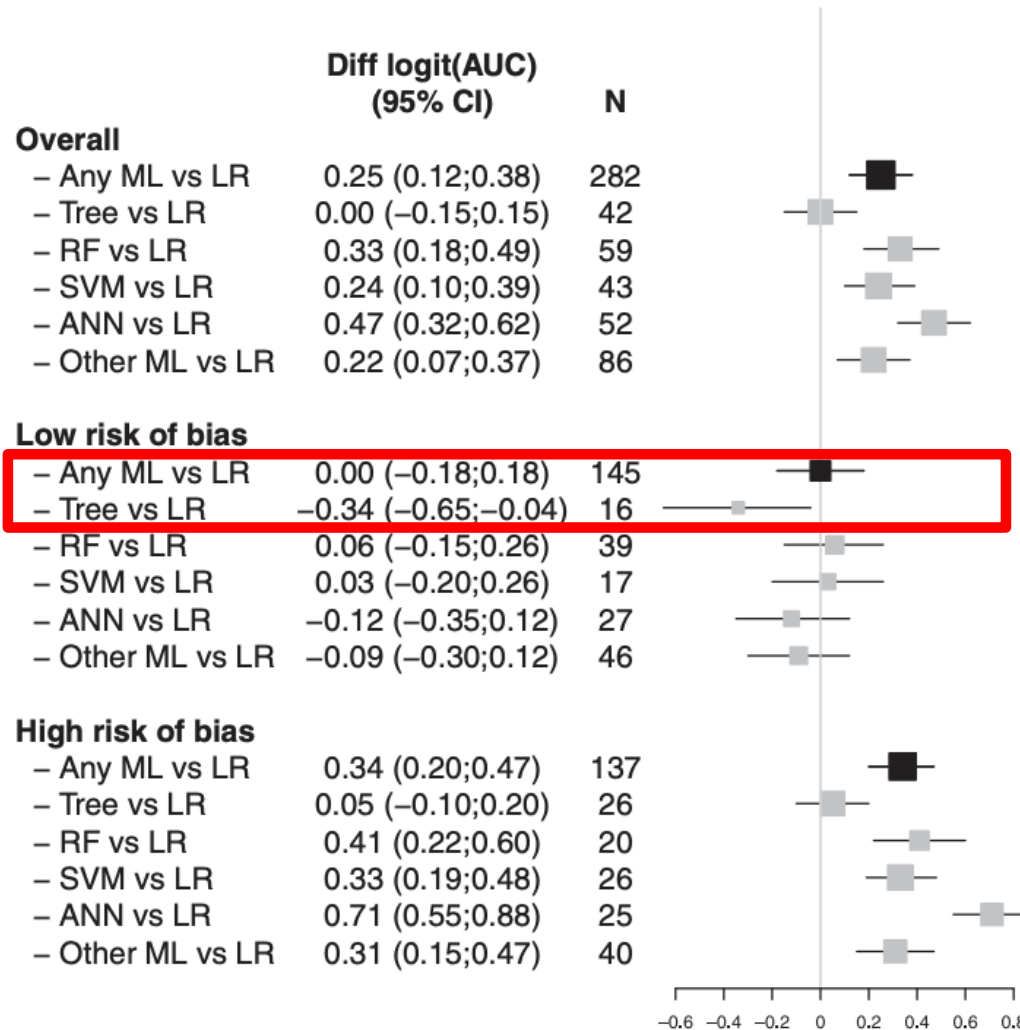
# Differences in discrimination



**Fig. 4.** Differences in discriminative ability between LR and ML models, overall and according to risk of bias ($n = 282$ comparisons).

**Arjun (Raj) Manrai**
@arjunmanrai

(Thread) The paper by Evangelia et al. in @JClinEpi on 'logistic regression = machine learning' for medicine has generated many reactions. This paper may be misinterpreted by #MachineLearning cynics and enthusiasts alike

**Arjun (Raj) Manrai** @arjunmanrai · 12 feb.

There are notable absences, such as many of the seminal contributions of deep learning to image analysis in medicine (e.g. Gulshan et al. JAMA 2016 and Esteva et al. Nature 2017). 7/n

Original Investigation | Innovations in ure
rnal of science

December 13, 2016

**Development and Vali**d: 25 January 2017
**Detection of Diabetic** tologist-level classification
**Photographs** with deep neural network

Varun Gulshan, PhD[1]; Lily Peng, MD, PhD[1]; Marc Coram
Brett Kuprel ✉, Roberto A. Novoa ✉, Justin Ko, Susan M. Swett
✉

≫ Author Affiliations | Article Information

*JAMA.* 2016;316(22):2402-2410. doi:10.1001/jama.201
118 (02 February 2017) | Download Citation ⬇

🌐 **Machine Learning Website**
um to this article was published on 28 June 2017

# Where is ML useful?

Large

A          B

N cases     Small ────────────────── Large

C          D

Small

N predictors

**Maarten van Smeden** @MaartenvSmeden · 27 jun.

Interesting correspondence about **machine learning** and the signal:noise ratio in @NEJM by **@BenVanCalster** @laure_wynants nejm.org/doi/full/10.10...

What do you think? The *advantage* of modern **machine learning** over traditional statistical approaches is more in....

| | |
|---|---|
| high signal:noise | 43% |
| low signal:noise | 32% |
| how dare you ask? | 25% |

178 stemmen · Eindresultaten

**Table 2. Key Questions to Ask When Deciding What Type of Model Is Necessary.**

**How complex is the prediction task?**

Simple prediction tasks are defined as those that can be performed with high accuracy with a small number of predictor variables. For example, predicting the development of hyperkalemia might be possible from just a small set of variables, such as renal function, the use of potassium supplements, and receipt of certain medications.

Complex prediction tasks are defined as those that cannot be predicted accurately with a small number of predictor variables. For example, identification of abnormalities in a pathological slide requires evaluation of patterns that are not obvious over millions of pixels.

In general, simple prediction tasks can be performed with traditional models (e.g., logistic regression), and complex tasks require more complex models (e.g., neural networks).

Rajkomar et al. NEJM 2019;380:1347-58.

# Myth 5: ML / AI leads to better generalizability

## Calibration drift in regression and machine learning models for acute kidney injury FREE

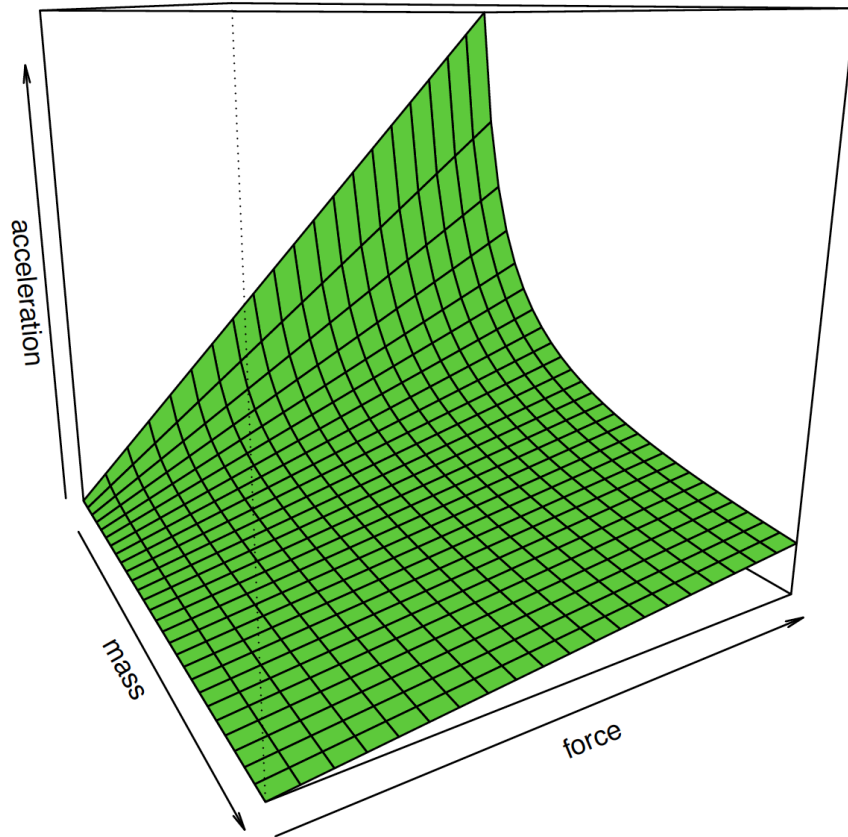Sharon E Davis, Thomas A Lasko, Guanhua Chen, Edward D Siew, Michael E Matheny ✉

" … developed 7 parallel models for hospital-acquired acute kidney injury using common regression and machine learning methods, validating each over 9 subsequent years.":
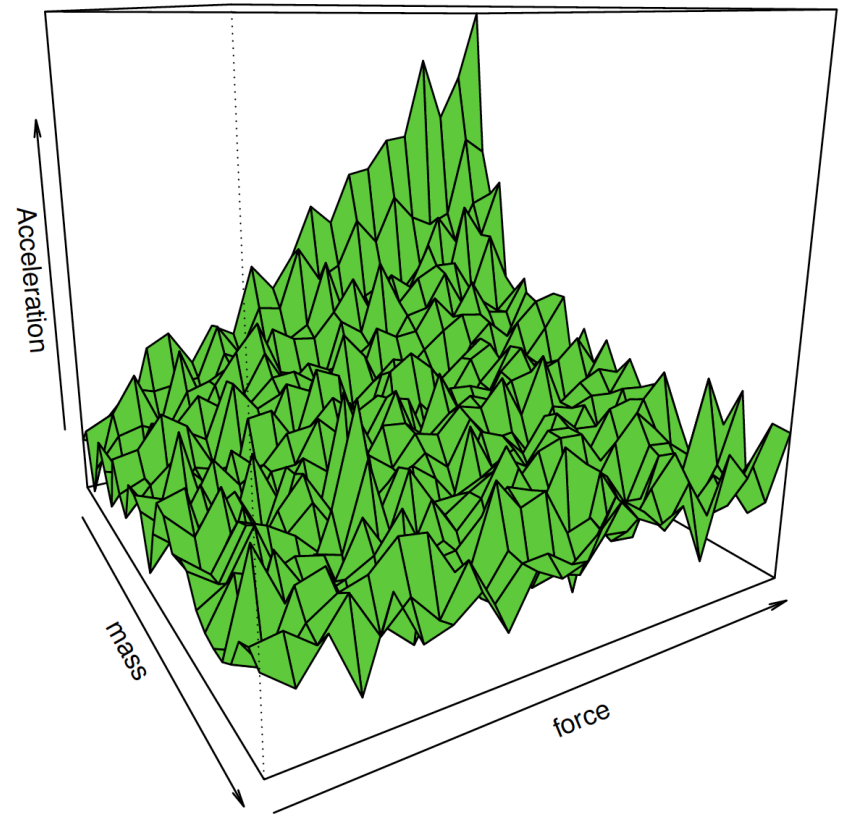
"Discrimination was maintained for all models. Calibration declined as all models increasingly overpredicted risk. **However, the random forest and neural network models maintained calibration** … "
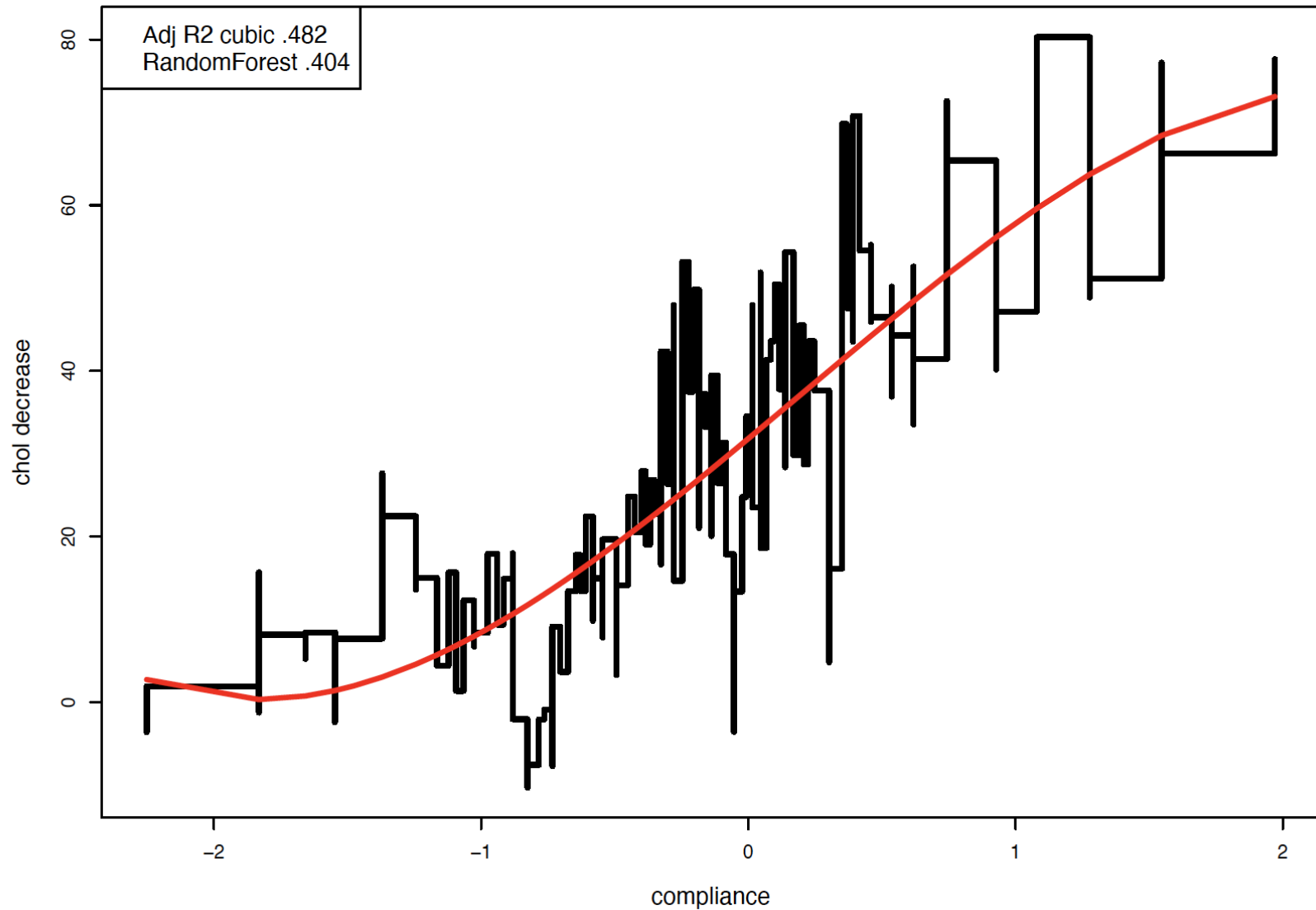
# Efron talk Leiden

**Newton's 2nd law: acceleration=force/mass**



**If Newton had done the experiment**

Cholesterol data: randomForest estimate (X=poly(c,8)), 500 trees,
compared with cubic regression curve

Adj R2 cubic .482
RandomForest .404

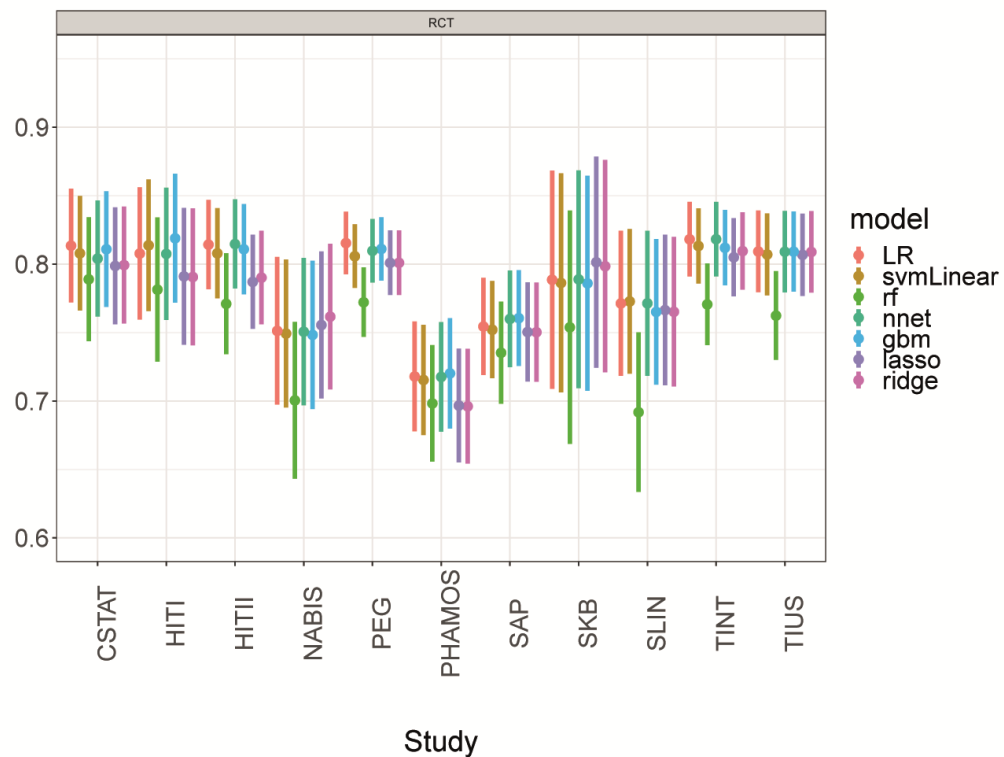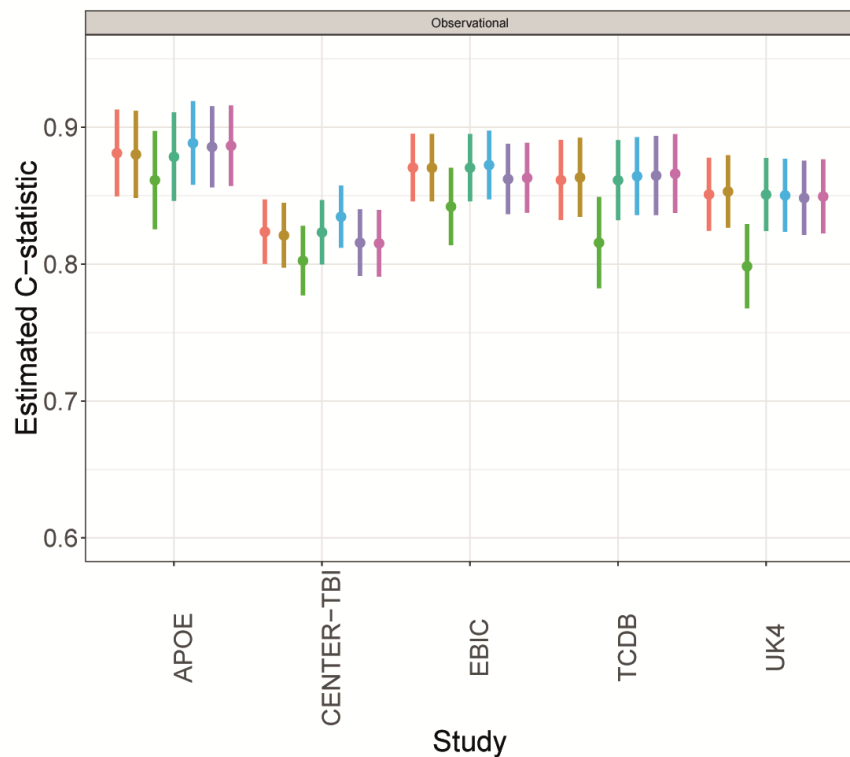chol decrease

compliance

# Empirical findings in TBI

– 16 cohorts: 5 observational, 11 RCTs

– Develop in 15, validate in 1

– 7 methods: LR; SVM; RF; nnet; gbm; LASSO; ridge

# 5 observational    11 RCTs



Variability between cohorts >> variability between methods

# Prediction challenges

- There is no such thing as a validated prediction algorithm

- Algorithms are high maintenance
  - Developed models need **validation and updating** to remain useful over time and place

- Regulation and quality control of algorithms
  - What about proprietary algorithms?

# Five myths

1. Big Data will resolve the problems of small data
   **NO: Big Data, Big Errors**

2. ML/AI is very different from classical modeling
   **NO: a continuum, cultural differences**

3. Deep learning is relevant for all medical prediction
   **NO: Deep learning excels in visual tasks**

4. ML / AI is better than classical modeling for prediction
   **NO: some methods do harm (e.g. tree modeling)**

5. ML / AI leads to better generalizability
   **NO: any prediction model may suffer from poor generalizability**