# Validation and Trustworthiness of AI based Predictions

**March 2024, Aachen**

Ewout W. Steyerberg, PhD

*Professor of Clinical Biostatistics and
Medical Decision Making*

Dept of Biomedical Data Sciences

Leiden University Medical Center

Thanks to many for assistance and inspiration

# "10% risk"

# Clinical prediction models

**EVIDENCIO**
MEDICAL DECISION SUPPORT

Models ⌄    Validations ⌄    About ⌄    Pricing ⌄

Find models by title, author, specialty, MeSH, ...

### Public models by Specialty

| | | |
|---|---|---|
| Adolescent medicine (5) | Immunology (4) | Paediatrics (6) |
| Aerospace medicine (0) | Infectious disease (11) | Palliative medicine (1) |
| Allergology (6) | Intensive care (32) | Pathology (2) |
| Anaesthesiology (2) | Internal medicine (44) | Physiatry (0) |
| Cardiology (76) | Microbiology (2) | Physical therapy (7) |
| Clinical chemistry (3) | Neonatology (4) | Podiatry (1) |
| Clinical pharmacology (22) | Nephrology (19) | Psychiatry (3) |
| Dermatology (1) | Neurology (14) | Psychotherapy (0) |
| Emergency medicine (32) | Neurophysiology (0) | Public Health (27) |
| Endocrinology (1) | Neuroradiology (0) | Pulmonology (32) |
| Epidemiology (7) | Neurosurgery (0) | Radiology (5) |
| Gastroenterology (20) | Nuclear medicine (2) | Radiotherapy (2) |
| General practice (59) | Obstetrics (8) | Rheumatology (4) |
| Geriatrics (37) | Occupational therapy (1) | Sports medicine (0) |
| Gerontology (1) | Oncology (182) | Surgery (96) |
| Gynaecology (23) | Ophthalmology (0) | Traumatology (18) |
| Health informatics (4) | Orthodontics (0) | Unspecified (46) |
| Hematology (6) | Orthopaedics (13) | Urology (63) |
| Hepatology (4) | Otorhinolaryngology (1) | Vascular medicine (11) |

# Validated → Trustworthy?

**3-year survival after resection in patients with pancreatic cancer: updated amsterdam model**

This prediction model for 3-year survival after resection in patients with pancreatic cancer was developed in 2015 in Amsterdam UMC (location AMC) and externally validated in 2019 in 3081 patients from 8 countries (USA, UK, Germany, Italy, Sweden, the Netherlands, Korea, Australia).PubMed ID 3192443...

Filtered by: ✕ CE Certified

There were no models found matching your search criteria

# Predictive algorithms: Medical AI



**Phase 1**

Data preparation

**Phase 2**

Development AI algorithm

**Phase 3**

Validation AI algorithm

**Phase 4**

Software environment

**Phase 5**

Impact assessment

**Phase 6**

Implementation in medical practice

## Guidelines and quality criteria for artificial intelligence-based prediction models in healthcare: a scoping review

Anne A. H. de Hond[1,2,3,8]✉, Artuur M. Leeuwenberg [4,8]✉, Lotty Hooft[4,5], Ilse M. J. Kant[1,2,3], Steven W. J. Nijman [4],

# Trustworthiness of predictions

Mathematical models needed:
- Complex processes
- No simple prediction via a deterministic theory

Modelling assumptions:
- Generally false
- Intelligent guesswork

Medical prediction: $y \sim X$

# Topics: trustworthy predictions

1. What do we need for individual patients?

- Internal vs external validity

- Calibration vs discrimination

2. Trustworthy processes to build a prediction model?

- AI

- Humans

- Requirements

3. Types of uncertainty

- Statistical aspects

- Model uncertainty

- Heterogeneity between contexts of practical application

# Validated = trustworthy?

- Classic:

  - No validation, only internal → low ranking journal

  - 1 convincing validation → top journal

- Modern

  - Substantial heterogeneity in performance

  → There is no such thing as a validated model

well calibrated          over prediction



Modification of `val.prob()` in rms; **`val.prob.ci.2()`**

*Steyerberg et al, PLoS Med 2008*

**Appendix 8: Calibration plot: observed proportion vs predicted probability of the clinical prediction model for 5 internal-external cross-validations.**
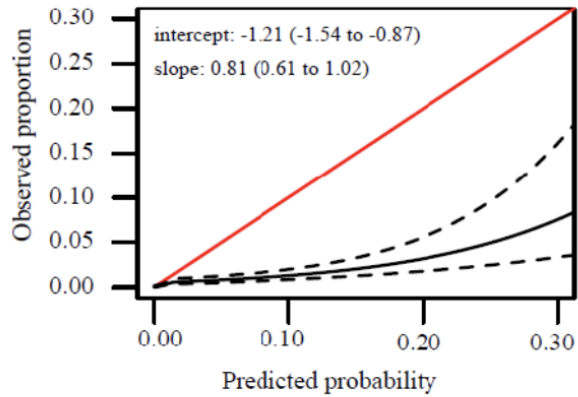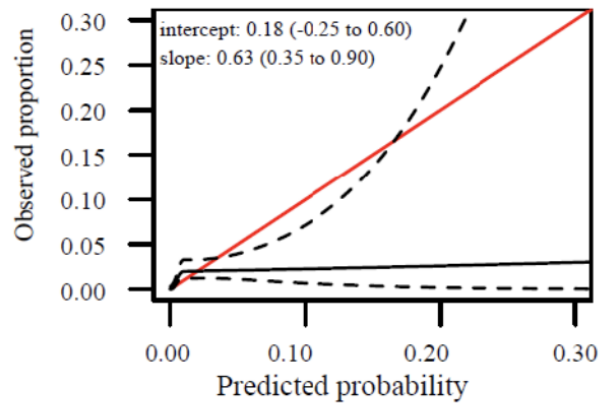
**COMMENT**   **OPEN**

# Perspectives on validation of clinical predictive algorithms

Anne A. H. de Hond [1,2,3 ✉], Vaibhavi B. Shah [2], Ilse M. J. Kant[4], Ben Van Calster [3,5], Ewout W. Steyerberg [1,3] and Tina Hernandez-Boussard [2,6,7]

**Fig. 1  Generalizability types.** Schematic overview of the different types of generalizability with the validation's goals, methods, and stakeholders.

# Summary validation

- Internal: minimum, same underlying population
- External:
    - Temporal
    - Geographic
    - Domain
- Efficient: internal-external cross-validation

# Performance assessment

- What is the most commonly reported measure for performance of prediction models?

    - Area under the Receiver Operating Characteristic curve (AUC), or concordance (c) statistic

    - Discrimination = spread of predictions between individuals

        - Higher if better predictors in model

        - Higher --> more trustworthy??

**Ben Van Calster** @BenVanCalster · 25 mei

I searched the literature some day about how people think about different AUC values. Interesting…

# Trustworthiness for individuals

- Calibration = reliability of predictions per individual

# well calibrated



CRASH GCS<=12 HIC, n=1351

Modification of `val.prob()` in rms; **val.prob.ci.2()**

*Steyerberg et al, PLoS Med 2008*

- Calibration = reliability of predictions per individual

  - True risk estimates UTOPIAN

## A calibration hierarchy for risk models was defined: from utopia to empirical data

Ben Van Calster[a,b,*], Daan Nieboer[b], Yvonne Vergouwe[b], Bavo De Cock[a], Michael J. Pencina[c,d], Ewout W. Steyerberg[b]

  - Calibration underreported

# Calibration: the Achilles heel of predictive analytics

Ben Van Calster[1,2,6*], David J. McLernon[3,6], Maarten van Smeden[2,4,6], Laure Wynants[1,5], Ewout W. Steyerberg[2,6]
On behalf of Topic Group 'Evaluating diagnostic tests and prediction models' of the STRATOS initiative[6]

Check for updates

# Topics: trustworthy predictions

1. What do we need for individual patients?

- Internal vs external validity

- Calibration vs discrimination

2. Trustworthy processes to build a prediction model?

- AI

- Humans

3. Types of uncertainty

- Statistical aspects

- Model uncertainty

- Heterogeneity between contexts of practical application

# Trustworthiness of ChatGPT

ChatGPT: may be hallucinating

- Simple calculations:



**can you do calculations?**

Of course! I can help with various types of calculations. What specific calculations would you like assistance with?

35623 × 472819

The product of 35,623 multiplied by 472,819 is:

$$35,623 \times 472,819 = 16,814,642,337.$$

35623 x 472819 = 16.843.231.237

**are you sure about this calculation?**

Apologies for the mistake. Let me correct that for you.

$$35,623 \times 472,819 = 16,814,506,737.$$

My apologies for the earlier response. Let me recalculate that for you:

35623 × 472819 = 16,827,062,037

So, the correct product of 35,623 multiplied by 472,819 is 16,827,062,037.

🖩 35623 x 472819 = 16.843.231.237

# Trustworthiness and AI

- Relation to evidence base
- Other popular terms related to AI
  - Fairness
  - Equity
- → Let's ask ChatGPT

what are important ethical concerns regarding AI for predictive purposes?

Ethical concerns regarding AI for predictive purposes are numerous and evolving as AI technologies advance. Here are some important ethical considerations:

1. **Bias and Fairness:**

   AI models can perpetuate or even amplify biases present in the data they are trained on. Ensuring fairness and addressing bias is a critical concern, especially in predictive models related to sensitive aspects like race, gender, or socioeconomic status.

2. **Privacy:**

   Predictive AI often relies on vast amounts of personal data. Safeguarding privacy, obtaining informed consent, and ensuring data is used appropriately and securely are essential ethical considerations.

3. **Transparency and Explainability:**

   Understanding how an AI model arrives at its predictions is crucial for trust and accountability. Models should be transparent and explainable to users, making it easier to identify potential biases and errors.

4. **Accountability and Responsibility:**

   Defining responsibility and accountability for the actions and decisions made by AI systems is challenging. Determining who is responsible when things go wrong, or when the AI system produces harmful predictions, is a key ethical concern.

5. **Autonomy and Decision-making:**

   Predictive AI can influence decision-making processes in various domains, from lending to criminal justice. Striking the right balance between automation and human intervention, preserving human autonomy, and avoiding undue influence is vital.

6. **Security:**

   Ensuring the security of AI systems is essential to prevent malicious use, data breaches, and other potential risks associated with predictive AI.

7. **Long-term Effects and Unintended Consequences:**

   Predictive AI can have unforeseen and far-reaching consequences, including social, economic, and cultural impacts. Anticipating and mitigating these consequences is a significant ethical consideration.

8. **Informed Consent and User Rights:**

   Obtaining meaningful informed consent from individuals whose data is used in predictive models is important. Individuals should have rights over their data and should be informed about how their data is being used.

9. **Data Quality and Representativeness:**

   Ensuring that the data used to train predictive models is of high quality and representative of the population it aims to predict for is essential to avoid skewed or inaccurate predictions.

10. **Human Augmentation vs. Replacement:**

    Balancing the use of AI for predictive purposes to enhance human capabilities (augmentation) versus replacing human decision-making entirely raises important ethical questions regarding the role of AI in society.

Addressing these ethical concerns is crucial to ensure the responsible and beneficial development and deployment of AI for predictive purposes.

**Can you group these concerns further to fewer groups?**

Certainly! We can group these ethical concerns into broader categories view:

1. **Data and Model Integrity:**

   a. **Bias and Fairness**

   b. **Data Quality and Representativeness**

   c. **Transparency and Explainability**

2. **Privacy and Consent:**

   a. **Privacy**

   b. **Informed Consent and User Rights**

3. **Accountability and Responsibility:**

   a. **Accountability for AI Actions**

   b. **Security**

   c. **Long-term Effects and Unintended Consequences**

4. **Human-Centric Concerns:**

   a. **Autonomy and Decision-making**

   b. **Human Augmentation vs. Replacement**

# Trustworthy models?

- Modeling flexibility: friend or foe?

  Human oversight on:

  - Classical modeling: selection of predictors; nonlinearity; interactions

  - AI: hyperparameters; technique CART / RF / XGBoost / nnet / ..

# Reproducibility trial: 246 biologists get different results from same data sets

**Wide distribution of findings shows how analytical choices drive conclusions.**

# Trustworthiness: poor for human modelers

Red cards and dark skin soccer players

https://psyarxiv.com/qkwst/

*Empirical Article*

ASSOCIATION FOR PSYCHOLOGICAL SCIENCE

## Many Analysts, One Data Set: Making Transparent How Variations in Analytic Choices Affect Results

R. Silberzahn[1], E. L. Uhlmann[2], D. P. Martin[3], P. Anselmi[4], F. Aust[5], E. Awtrey[6], Š. Bahník[7], F. Bai[8], C. Bannard[9], E. Bonnier[10], R. Carlsson[11], F. Cheung[12], G. Christensen[13], R. Clay[14], M. A. Craig[15], A. Dalla Rosa[4], L. Dam[16], M. H. Evans[17], I. Flores Cervantes[18], N. Fong[19], M. Gamez-Djokic[20], A. Glenz[21], S. Gordon-McKeon[22], T. J. Heaton[23], K. Hederos[24], M. Heene[25], A. J. Hofelich Mohr[26], F. Högden[5], K. Hui[27], M. Johannesson[10], J. Kalodimos[28], E. Kaszubowski[29], D. M. Kennedy[30], R. Lei[15], T. A. Lindsay[26], S. Liverani[31], C. R. Madan[32], D. Molden[33], E. Molleman[16], R. D. Morey[34], L. B. Mulder[16], B. R. Nijstad[16], N. G. Pope[35], B. Pope[36], J. M. Prenoveau[37], F. Rink[16], E. Robusto[4], H. Roderique[38], A. Sandberg[24], E. Schlüter[39], F. D. Schönbrodt[25], M. F. Sherman[37], S. A. Sommer[40], K. Sotak[41], S. Spain[42], C. Spörlein[43], T. Stafford[44], L. Stefanutti[4], S. Tauber[16], J. Ullrich[21], M. Vianello[4], E.-J. Wagenmakers[45], M. Witkowiak[46], S. Yoon[19], and B. A. Nosek[3,47]

[1]Organisational Behaviour, University of Sussex Business School; [2]Organisational Behaviour Area, INSEAD Asia Campus;

- 29 teams involving 61 analysts; same dataset; same research question: **whether soccer referees are more likely to give red cards to dark skin toned players than light skin toned players**
- Estimated odds ratios 0.89 –2.93 (median 1.3)
- 20 teams: statistically significant positive effect, 9: non-significant relation

# Estimated odds ratios by 29 research teams

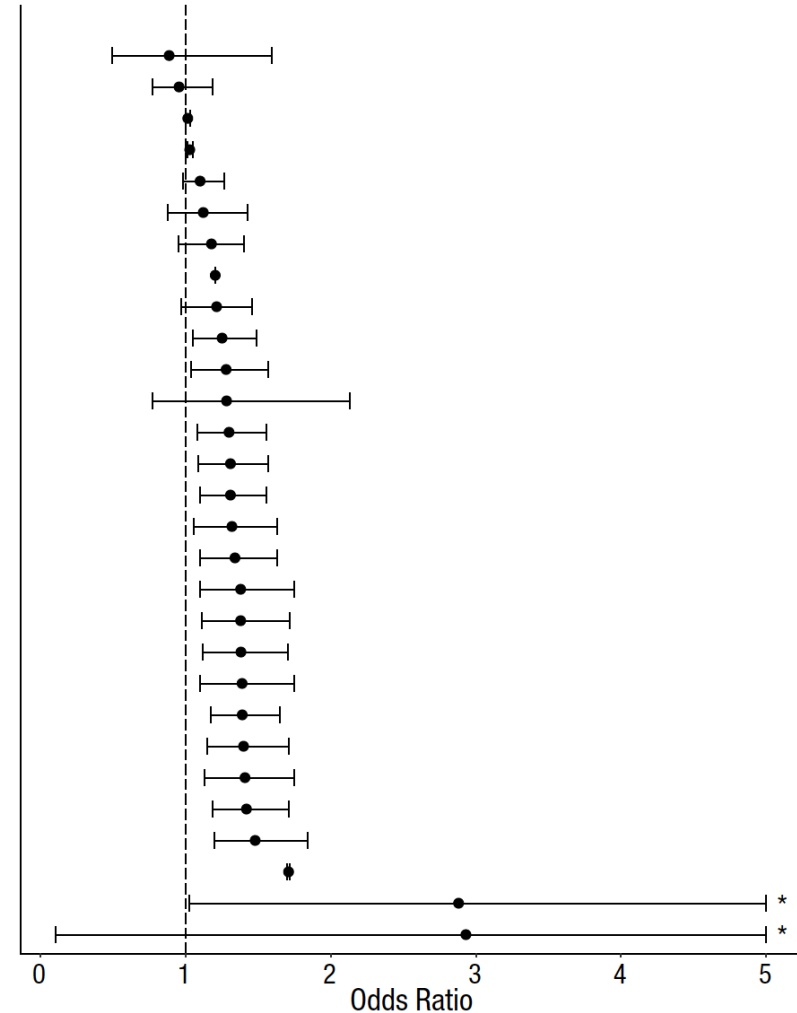| Team | Analytic Approach | Odds Ratio |
|------|-------------------|------------|
| 12 | Zero-Inflated Poisson Regression | 0.89 |
| 17 | Bayesian Logistic Regression | 0.96 |
| 15 | Hierarchical Log-Linear Modeling | 1.02 |
| 10 | Multilevel Regression and Logistic Regression | 1.03 |
| 18 | Hierarchical Bayes Model | 1.10 |
| 31 | Logistic Regression | 1.12 |
| 1 | OLS Regression With Robust Standard Errors, Logistic Regression | 1.18 |
| 4 | Spearman Correlation | 1.21 |
| 14 | WLS Regression With Clustered Standard Errors | 1.21 |
| 11 | Multiple Linear Regression | 1.25 |
| 30 | Clustered Robust Binomial Logistic Regression | 1.28 |
| 6 | Linear Probability Model | 1.28 |
| 26 | Hierarchical Generalized Linear Modeling With Poisson Sampling | 1.30 |
| 3 | Multilevel Logistic Regression Using Bayesian Inference | 1.31 |
| 23 | Mixed-Model Logistic Regression | 1.31 |
| 16 | Hierarchical Poisson Regression | 1.32 |
| 2 | Linear Probability Model, Logistic Regression | 1.34 |
| 5 | Generalized Linear Mixed Models | 1.38 |
| 24 | Multilevel Logistic Regression | 1.38 |
| 28 | Mixed-Effects Logistic Regression | 1.38 |
| 32 | Generalized Linear Models for Binary Data | 1.39 |
| 8 | Negative Binomial Regression With a Log Link | 1.39 |
| 20 | Cross-Classified Multilevel Negative Binomial Model | 1.40 |
| 13 | Poisson Multilevel Modeling | 1.41 |
| 25 | Multilevel Logistic Binomial Regression | 1.42 |
| 9 | Generalized Linear Mixed-Effects Models With a Logit Link | 1.48 |
| 7 | Dirichlet-Process Bayesian Clustering | 1.71 |
| 21 | Tobit Regression | 2.88 |
| 27 | Poisson Regression | 2.93 |



**Fig. 2.** Point estimates (in order of magnitude) and 95% confidence intervals for the effect of soccer players' skin tone on the number of red cards awarded by referees. Reported results, along with the analytic approach taken, are shown for each of the 29 analytic teams. The teams are ordered so that the smallest reported effect size is at the top and the largest is at the bottom. The asterisks indicate upper bounds that have been truncated to increase the interpretability of the plot; the actual upper bounds of the confidence intervals were 11.47 for Team 21 and 78.66 for Team 27. OLS = ordinary least squares; WLS = weighted least squares.
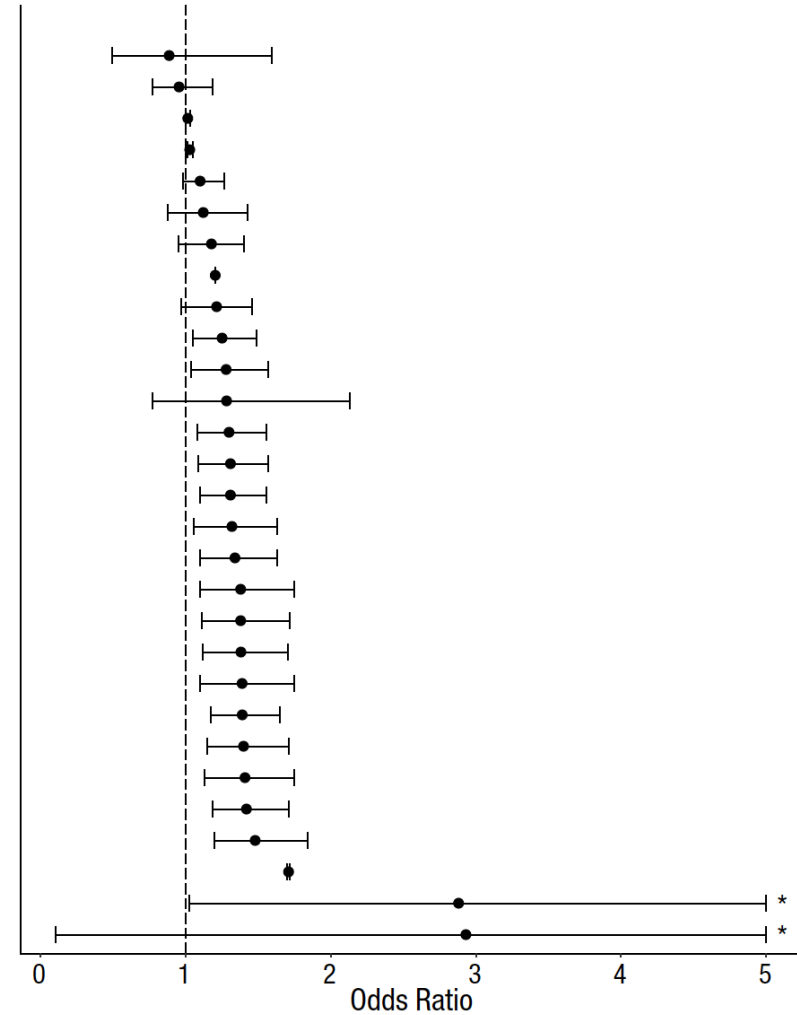
Fig. 2. Point estimates (in order of magnitude) and 95% confidence intervals for the effect of soccer players' skin tone on the number of red cards awarded by referees. Reported results, along with the analytic approach taken, are shown for each of the 29 analytic teams. The teams are ordered so that the smallest reported effect size is at the top and the largest is at the bottom. The asterisks indicate upper bounds that have been truncated to increase the interpretability of the plot; the actual upper bounds of the confidence intervals were 11.47 for Team 21 and 78.66 for Team 27. OLS = ordinary least squares; WLS = weighted least squares.

**Fig. 2.** Point estimates (in order of magnitude) and 95% confidence intervals for the effect of soccer players' skin tone on the number of red cards awarded by referees. Reported results, along with the analytic approach taken, are shown for each of the 29 analytic teams. The teams are ordered so that the smallest reported effect size is at the top and the largest is at the bottom. The asterisks indicate upper bounds that have been truncated to increase the interpretability of the plot; the actual upper bounds of the confidence intervals were 11.47 for Team 21 and 78.66 for Team 27. OLS = ordinary least squares; WLS = weighted least squares.

# Trustworthiness: poor for human modelers

- 29 teams involving 61 analysts; same dataset; same research question: whether soccer referees are more likely to give red cards to dark skin toned players than light skin toned players

- Estimated odds ratios 0.89 –2.93 (median 1.3).

- 20 teams: statistically significant positive effect, 9: non-significant relation.

- **21 unique combinations of covariates**

- **"Variation in analysis of complex data may be difficult to avoid, even by experts with honest intentions"**

*Corrigendum*

**ASSOCIATION FOR PSYCHOLOGICAL SCIENCE**

# Corrigendum: Many Analysts, One Data Set: Making Transparent How Variations in Analytic Choices Affect Results

" .. the authors forgot to add a citation of the *Nature* commentary to the final published version of the *AMPPS* article or to note that the main findings had been previously publicized via the commentary, the online preprint, research presentations at conferences and universities, and media reports by other people. The authors regret the oversight."

# Explainable to humans = trustworthy?

- Explainable AI

  - Algorithm trustworthy: if predictions are based on factors that are acceptable to domain experts instead of on 'spurious correlations'

  - SHAP (SHapley Additive exPlanations) values
    "By using SHAP values, researchers and practitioners can gain a deeper understanding of how different features influence model predictions, leading to improved model interpretability and trust."
    ChatGPT3.5

# Topics: trustworthy predictions

1. What do we need for individual patients?

- Internal vs external validity

- Calibration vs discrimination

2. Trustworthy processes to build a prediction model?

- AI

- Humans

3. Types of uncertainty

- Statistical aspects

- Model uncertainty

- Heterogeneity between contexts of practical application

# Approaches to uncertainty quantification

- Sample size

  - specifically **#events** for binary outcome prediction

- 'patients like you' and exceptionality

  - For risk communication: aleatory uncertainty

  - For uncertainty communication: epistemic

## COUNSELING MEN WITH PROSTATE CANCER: A NOMOGRAM FOR PREDICTING THE PRESENCE OF SMALL, MODERATELY DIFFERENTIATED, CONFINED TUMORS

MICHAEL W. KATTAN, JAMES A. EASTHAM, THOMAS M. WHEELER, NORIO MARU,

## PREDICTION OF INDOLENT PROSTATE CANCER

**Points**
0    10    20    30    40    50    60    70    80    90    100

**Pre.Tx.PSA**
20  13  9  7   5   4   3   2   1   0.5   0.2

**Pri.Bx.GI**
2
3

**Sec.Bx.GI**
2
3

**Total Points**
0    10    20    30    40    50    60    70    80    90    100    110

**Prob. Indolent Ca.**
0.08   0.12    0.2    0.3    0.4    0.5    0.6    0.7    0.8    0.9  0.92

**Instructions for Physician:** Locate the patient's PSA on the **PreTx PSA** axis. Draw a line straight upwards to the **Points** axis to determine how many points towards having an indolent cancer the patient receives for his PSA. Repeat this process for the remaining axes, each time drawing straight upward to the **Points** axis. Sum the points achieved for each predictor and locate this sum on the **Total Points** axis. Draw a line straight down to find the patient's probability of having indolent cancer.

**Instruction to Patient:** "Mr. X, if we had 100 men exactly like you, we would expect <predicted probability from nomogram * 100 > to have indolent cancer.
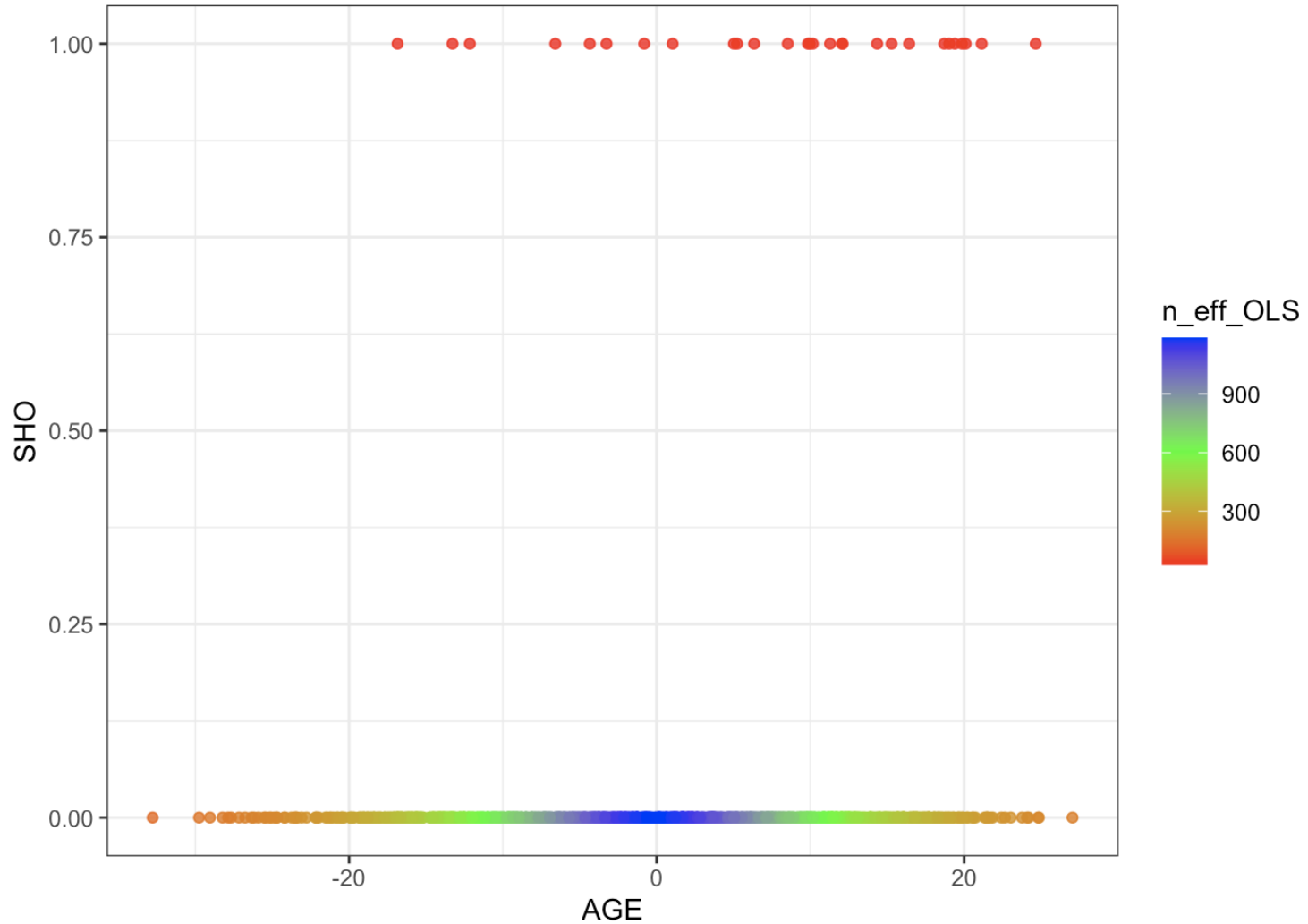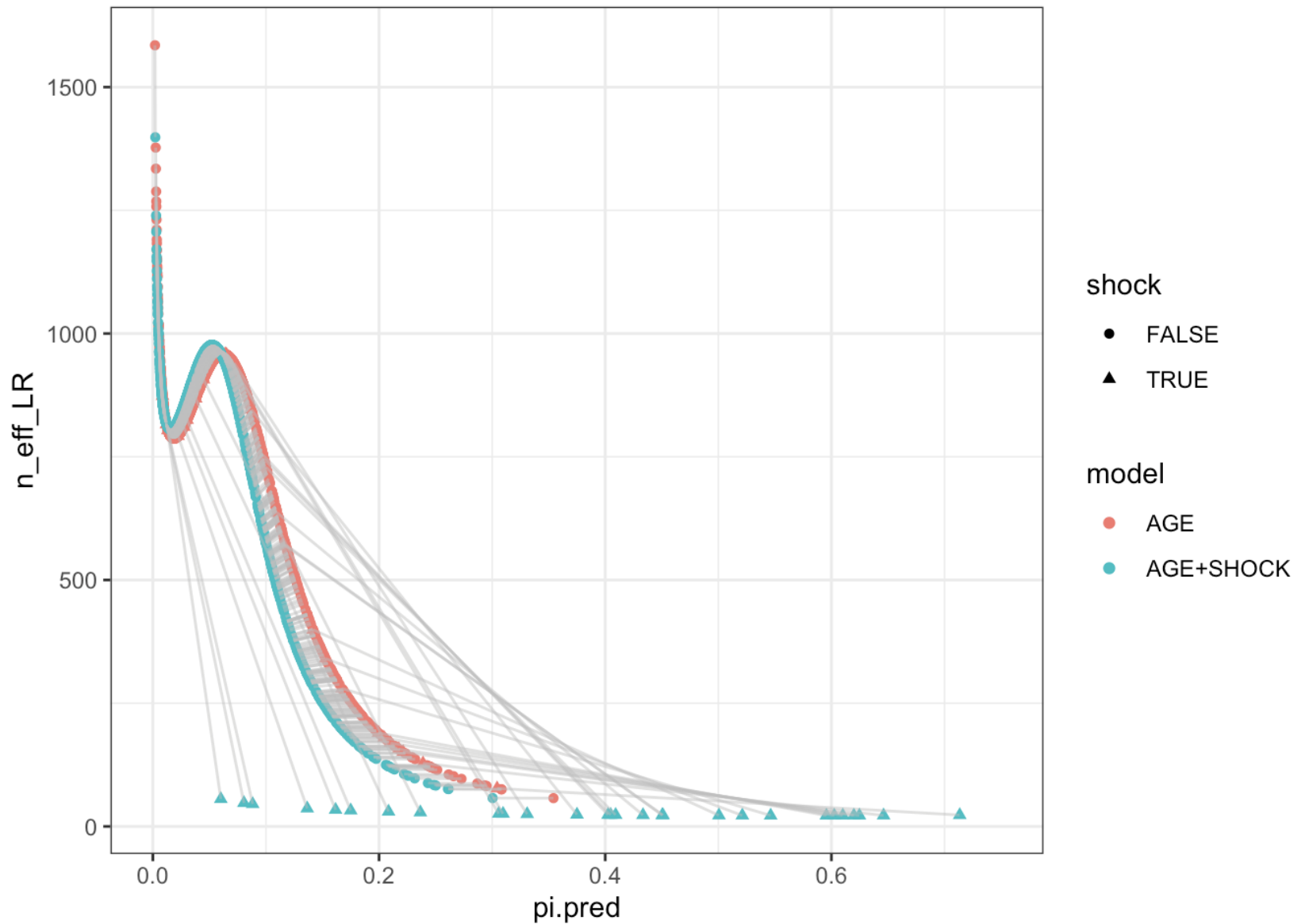
Effective n visualisations

Doranne Thomassen

2021-10-27

GUSTO data, n=1200
(out of 40,830)

# AI: arXiv paper

Kim de Bie [1]  Ana Lucic [1]  Hinda Haned [1]

"Trustworthiness expresses whether a prediction is aligned with the train data"

- Distance between the new patient and similar patients from the training data estimates the trustworthiness of a prediction; resembles the reference data



**To Trust or Not to Trust a Regressor**

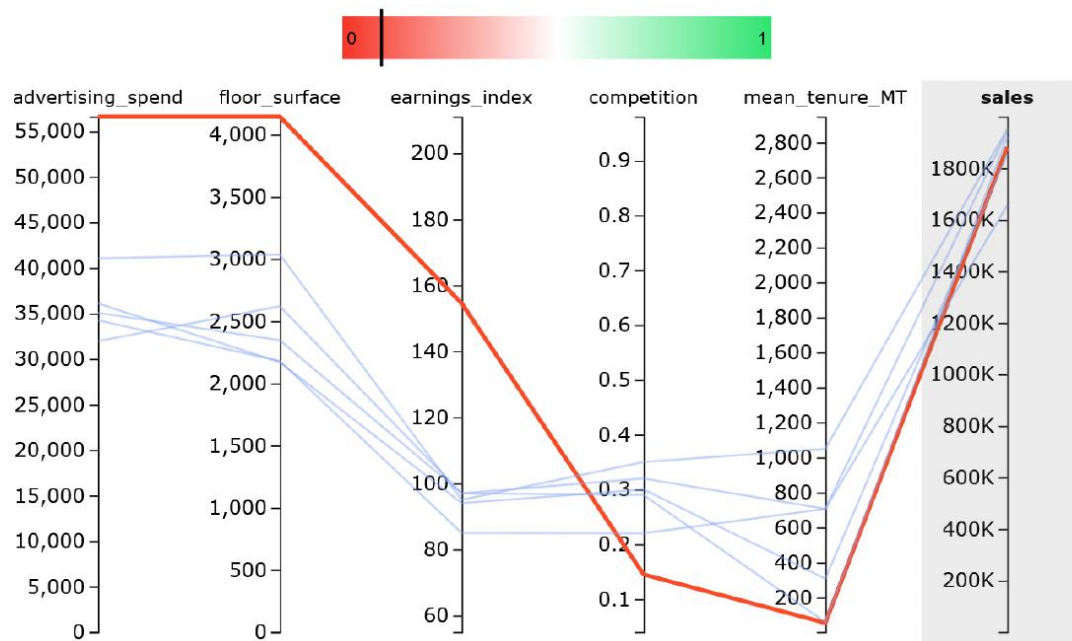The RETRO-score for this prediction is 0.091.

Figure 1. Example RETRO-VIZ output for an untrustworthy prediction. This model predicts sales based on five features. The RETRO-
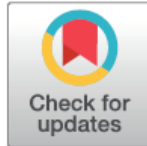
# Claim

"Effective N is an attractive concept to address epistemiologic uncertainty"

- Analytic solutions for regression models
    - Minimum certainty, say, n>10, for model specification: selection / shrinkage?
- Approximate solutions for machine learning models
    - Bootstrap

- Effective N: conditional on the model

# Model uncertainty

- 246 biologists modeling

- 61 analysts in 29 teams on the Red Card problem

- …

- Comparisons between classic vs machine learning

## 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,*]

**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
journal of science

December 13, 2016                    d: 25 January 2017

**Development and Vali** tologist-level classification
**Detection of Diabetic** with deep neural network
**Photographs**

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    118 (02 February 2017)  |  Download Citation ⬇

JAMA. 2016;316(22):2402-2410. doi:10.1001/jama.201

🌐 **Machine Learning Website**    um to this article was published on 28 June 2017
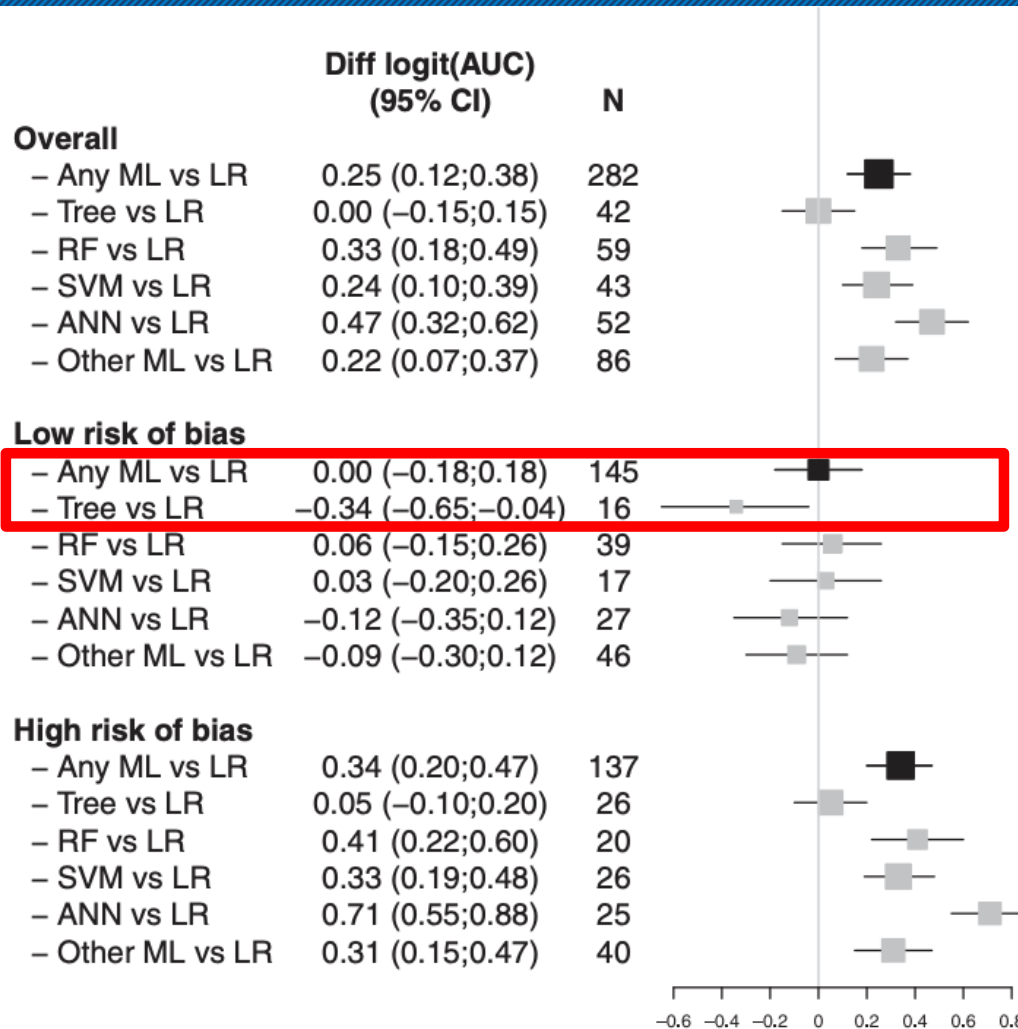
# Differences in discrimination



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

# Topics: trustworthy predictions

1. What do we need for individual patients?

- Internal vs external validity
- Calibration vs discrimination

2. Trustworthy processes to build a prediction model?

- AI
- Humans

3. Types of uncertainty

- Statistical aspects
- Model uncertainty
- Heterogeneity between contexts of practical application

# Heterogeneity

- Study design
- Selection of subjects
- Disease domain
- Measurement of covariates
- Measurement of outcomes
- Associations of covariates with outcome
- Overall outcome rates

# Heterogeneity in performance

| Performance measure | IMPACT models | | | |
|---|---|---|---|---|
| Mortality | Discrimination: 56 validations<br><br>Calibration: 31 validations | | | |
| | Mean[a] | | | Range |
| AUC | 0.79 | | | 0.65-0.90 |
| Calibration slope | 1.1 | | | 0.42-2.3 |
| Calibration intercept | -0.22 | | | -3.3-0.93 |

Steyerberg et al, *PLoS Med 2008*

Dijkland S et al; *J Neurotrauma 2019*

**RESEARCH ARTICLE**

WILEY Statistics in Medicine

# Assessment of heterogeneity in an individual participant data meta-analysis of prediction models: An overview and illustration

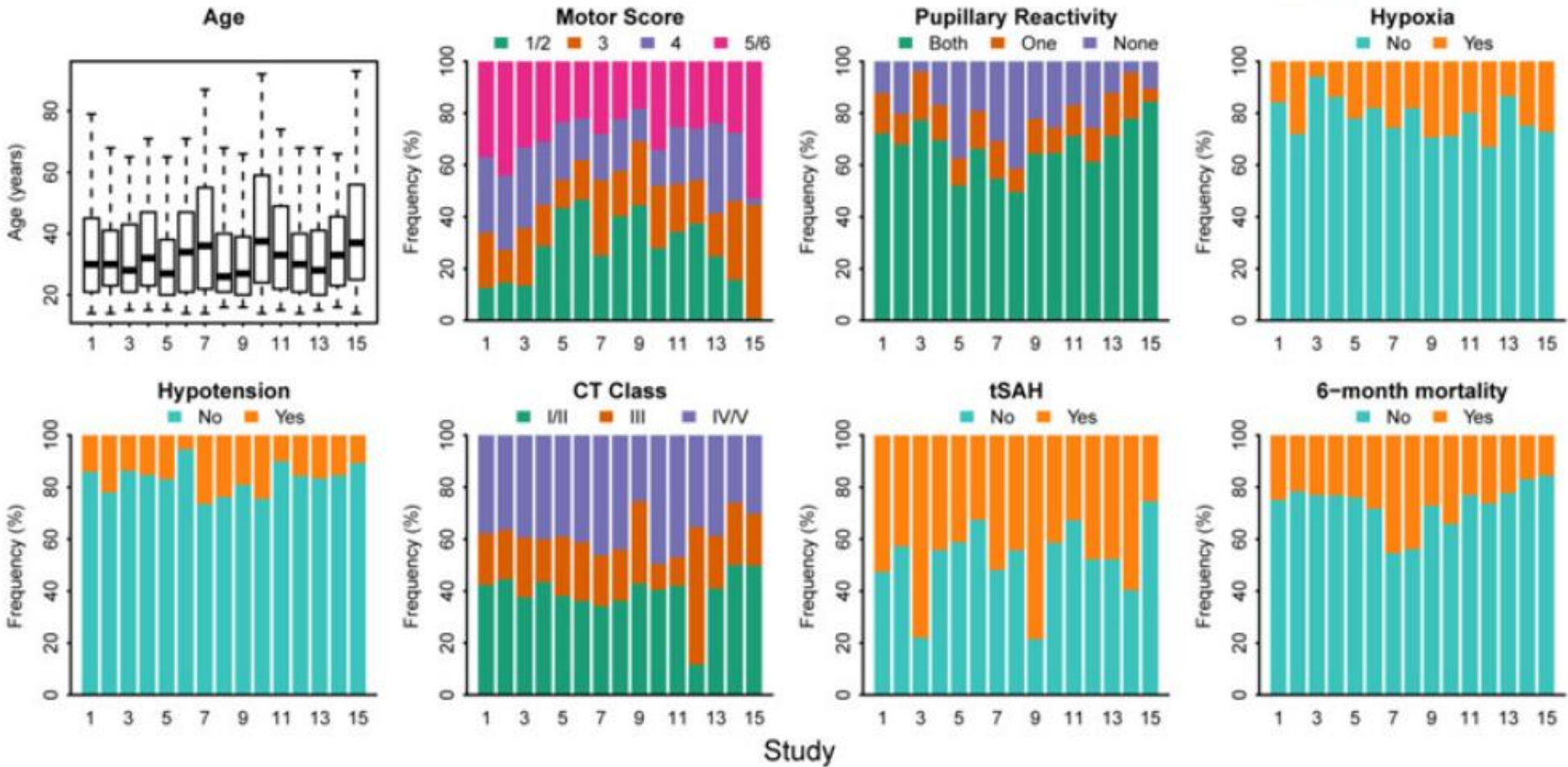Ewout W. Steyerberg[1,2] | Daan Nieboer[2] | Thomas P.A. Debray[3,4] |
Hans C. van Houwelingen[1]

# 15 cohorts: 11 RCTs, 4 Observational studies

| Nr. | Name | Enrollment period | Type[1] | n |
|---|---|---|---|---|
| 1 | TINT | 1991–1994 | RCT | 1118 |
| 2 | TIUS | 1991–1994 | RCT | 1041 |
| 3 | SLIN | 1994–1996 | RCT | 409 |
| 4 | SAP | 1995–1997 | RCT | 919 |
| 5 | PEG | 1993–1995 | RCT | 1510 |
| 6 | HIT I | 1987–1989 | RCT | 350 |
| 7 | UK4 | 1986–1988 | OBS | 791 |
| 8 | TCDB | 1984–1987 | OBS | 603 |
| 9 | SKB | 1996–1996 | RCT | 126 |
| 10 | EBIC | 1995–1995 | OBS | 822 |
| 11 | HIT II | 1989–1991 | RCT | 819 |
| 12 | NABIS | 1994–1998 | RCT | 385 |
| 13 | CSTAT | 1996–1997 | RCT | 517 |
| 14 | PHARMOS | 2001–2004 | RCT | 856 |
| 15 | APOE | 1996–1999 | OBS | 756 |

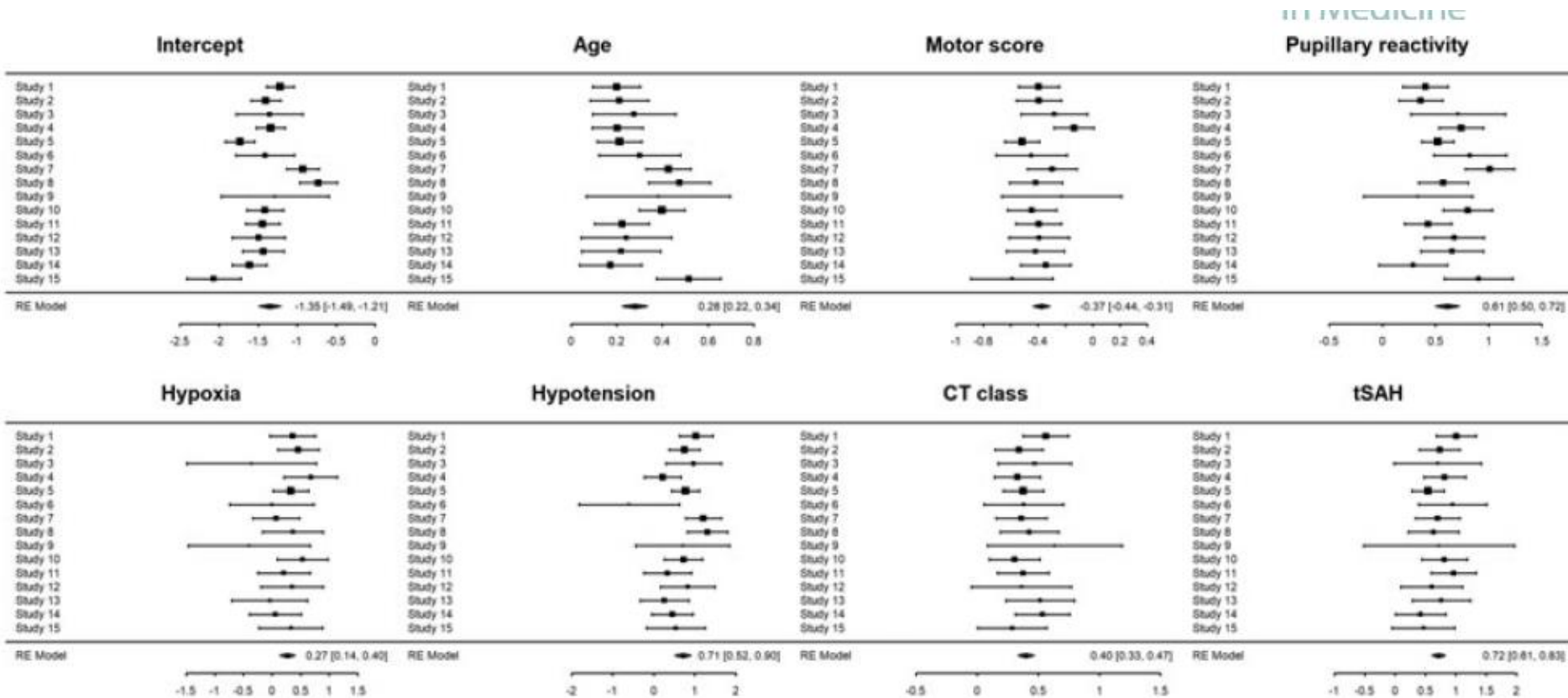[1]Type of study, RCT: randomized controlled trial, OBS: observational cohort

FIGURE 2    Forest plots showing estimated multivariable logistic regression coefficients and associated 95% confidence interval per study. The largest heterogeneity was noted for pupillary reactivity ($\tau = 0.17$) and hypotension ($\tau = 0.27$)
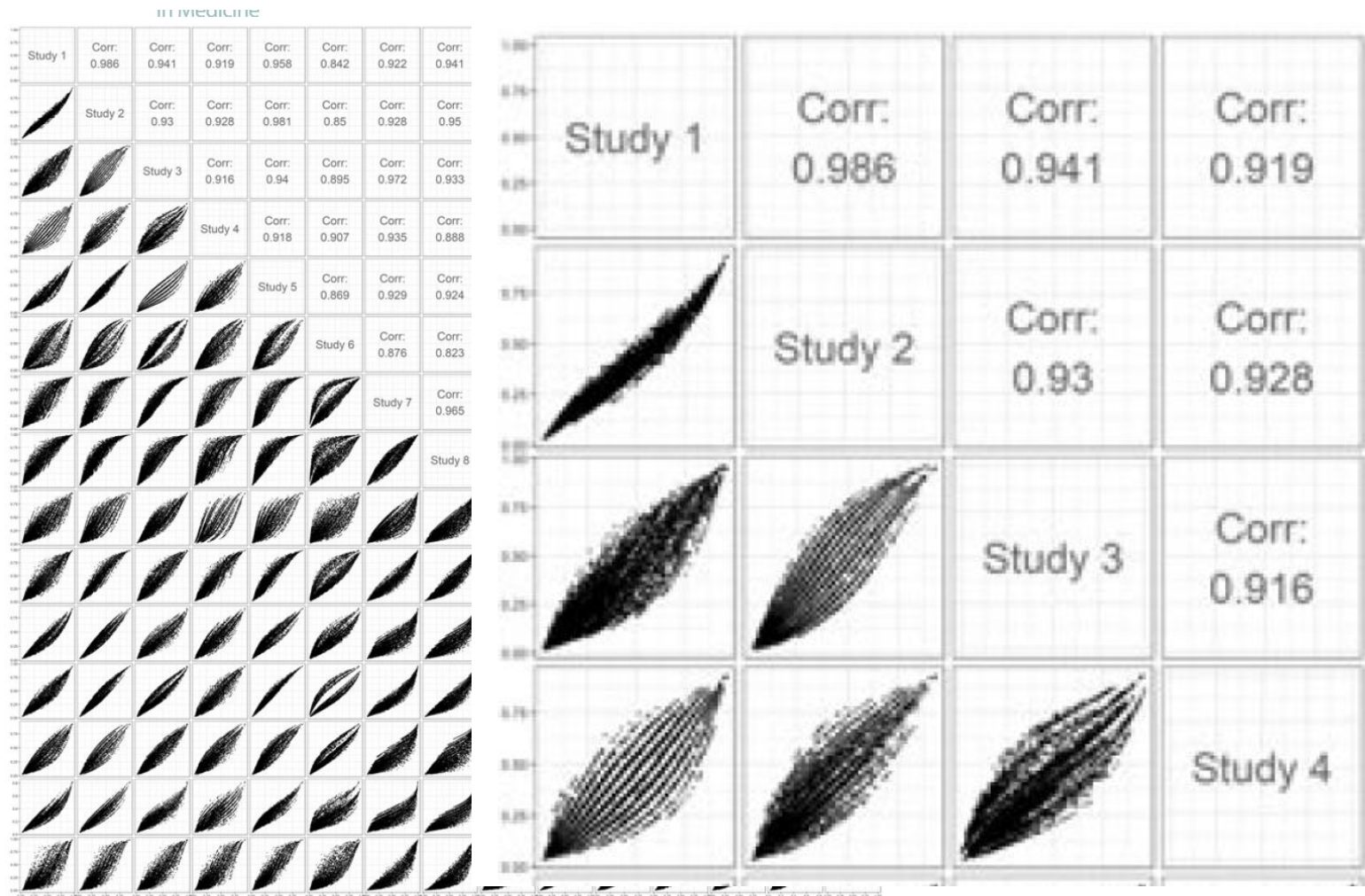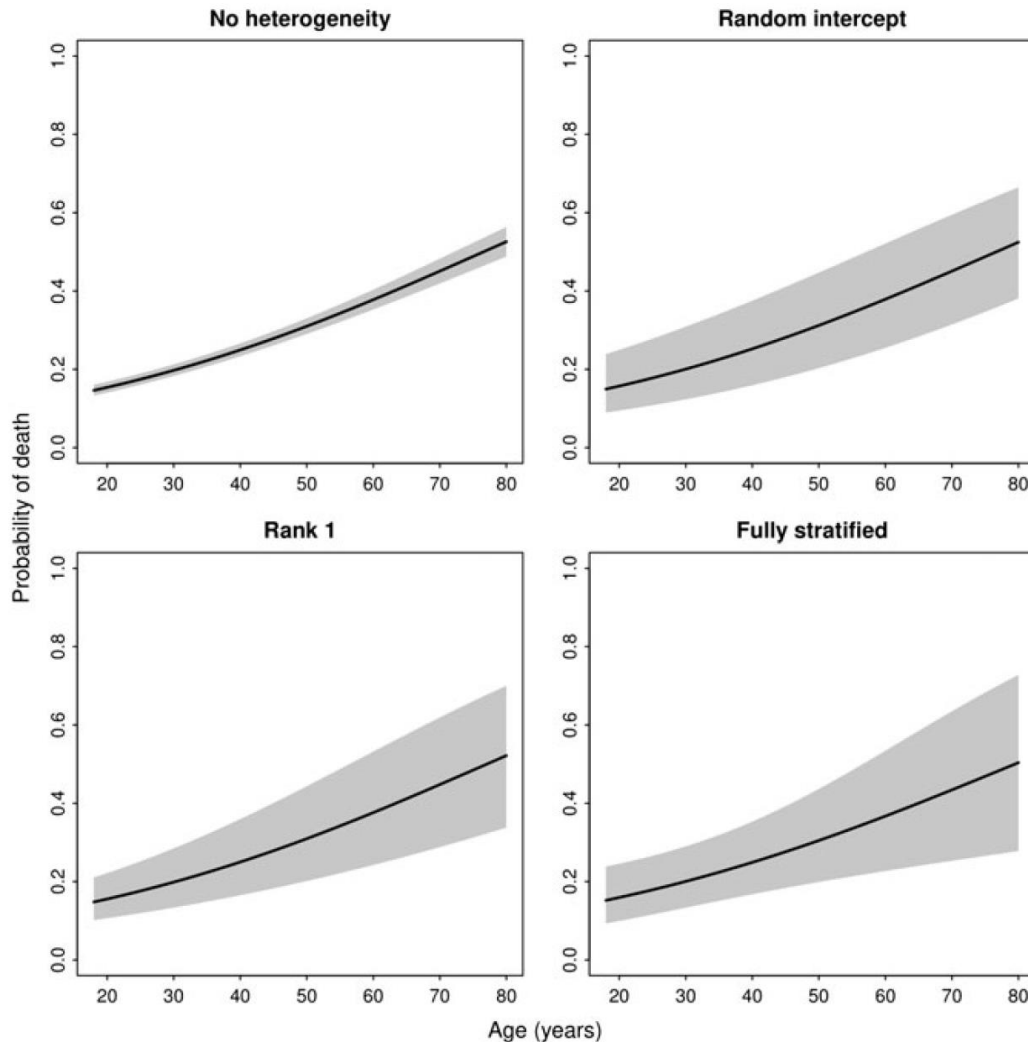
FIGURE 3 Correlation between predictions of study-specific models in a pairwise comparison between studies: 1-to-1 comparisons of predictions for all patients in the individual patient data set ($n = 11\,022$)

$$\text{Odds}(Y = 1 | \boldsymbol{X} = \boldsymbol{x}) = \exp(\alpha)\exp(\boldsymbol{X}\boldsymbol{\beta}').$$

$$\text{logit}\left(p_{ij}\right) = \alpha_j + \boldsymbol{X}_{ij}\boldsymbol{\beta}',$$
$$\alpha_j \sim N\left(\alpha, \tau_\alpha^2\right).$$

$$\text{logit}\left(p_{ij}\right) = \alpha_j + \gamma_j \boldsymbol{X}_{ij}\boldsymbol{\beta}',$$
$$\left(\alpha_j, \gamma_j\right) \sim MVN((\alpha, 1), \mathbf{T})$$

$$\text{logit}\left(p_{ij}\right) = \alpha_j + \boldsymbol{X}_{ij}\boldsymbol{\beta}_j',$$
$$\left(\alpha_j, \boldsymbol{\beta}_j\right) \sim MVN(\boldsymbol{\mu}, \mathbf{T}).$$

# Conclusions on trustworthy predictions

- Epistemic uncertainty: under the influence of the modeler

  - Larger sample sizes

  - Modest modeling, limit flexibility

- Heterogeneity: assess differences between settings

  - Study design

  - Distribution and effects of covariates

  - Differences between predictions

- Model predictions suffer from multiple sources of uncertainty

  - Transparency: for policy makers / physicians / patients

  - Context dependency: Local versus global models