

Assessment of performance and decision curve analysis

Ewout Steyerberg, Andrew Vickers

Dept of Public Health, Erasmus MC, Rotterdam, the Netherlands

Dept of Epidemiology and Biostatistics, Memorial Sloan-Kettering Cancer Center, New York, USA

Freiburg, October 2008

Work from Freiburg

BIOINFORMATICS

Assessment of survival prediction models based on microarray data

Martin Schumacher^{a *}, Harald Binder^b, Thomas Gerds^b

^aDepartment of Medical Biometry and Statistics, Institute of Medical Biometry and Medical Informatics, University Medical Center Freiburg, Germany

^b Freiburg Center of Data Analysis and Model Building, University Freiburg, Germany

WILEY

Multivariable Model-building

A pragmatic approach to regression
analysis based on fractional polynomials
for modelling continuous variables



Patrick Royston and Willi Sauerbrei

WILEY SERIES IN PROBABILITY AND STATISTICS

Erasmus MC

Erasmus

Erasmus MC – University Medical Center Rotterdam



Some issues in performance assessment

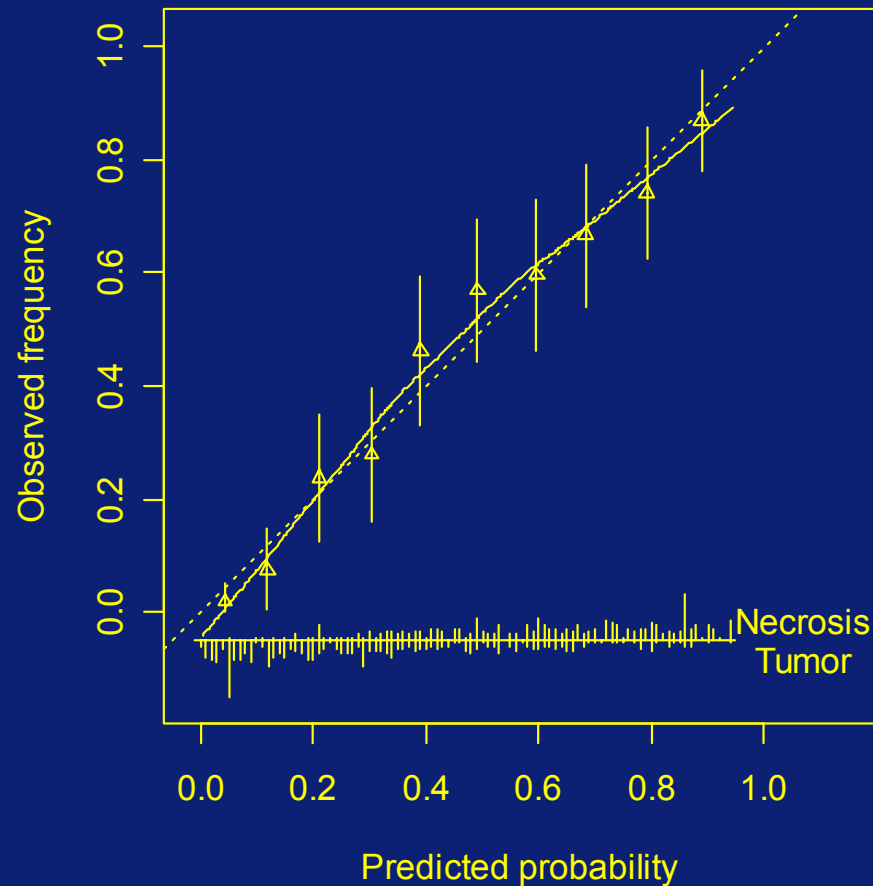
- Usefulness / Clinical utility: what do we mean exactly?
 - Evaluation of predictions
 - Evaluation of decisions
- Usefulness of a marker
 - Challenges in design and analysis
 - Measurement worth the increase in complexity (physician burden) and worth the costs (patient burden, financial costs)?
 - Additional value to a model with free / easy to obtain predictors
 - Validity of the model w/o marker

Traditional performance evaluation of predictions

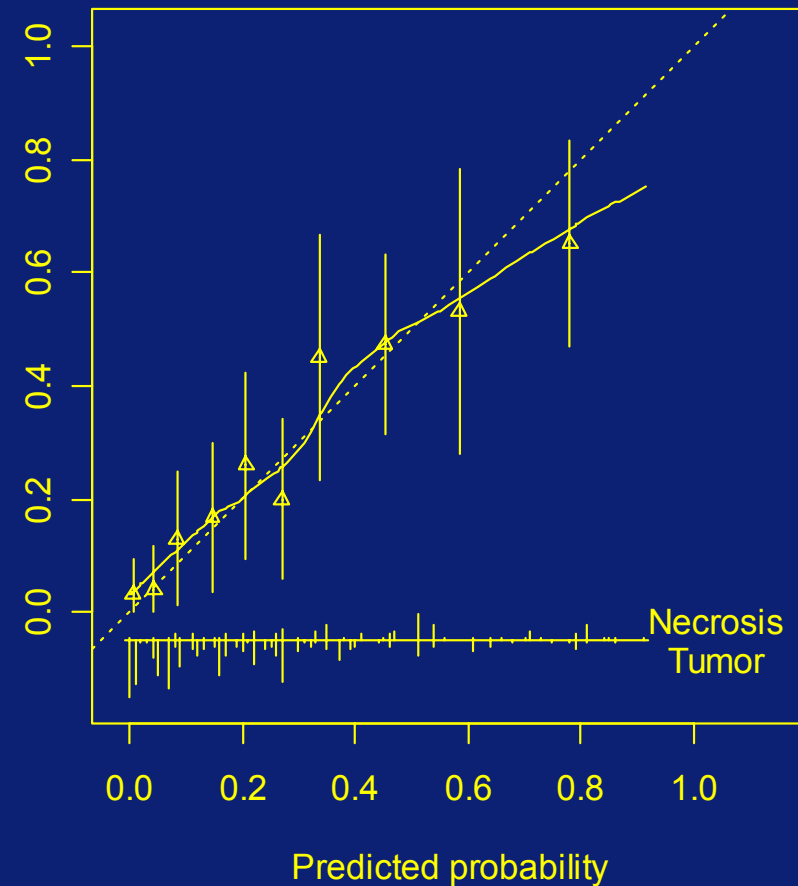
- Predictions close to observed outcomes?
 - Overall; consider residuals $y - \hat{y}$, or $y - p$
 - Brier score
 - R^2 (e.g. on log likelihood scale)
 - Discrimination: separate low risk from high risk
 - Area under ROC (or c statistic)
 - Calibration: e.g. 70% predicted = 70% observed
 - Calibration-in-the-large
 - Calibration slope

Validation graph to visualize both calibration, discrimination, and usefulness

Development, n=544



Validation, n=273



Statistics for Biology and Health

Ewout W. Steyerberg

Clinical Prediction Models

A Practical Approach to
Development, Validation, and
Updating

 Springer

Erasmus MC


Quantification of performance .. many developments

Brier score for model performance

BIOINFORMATICS

Assessment of survival prediction models based on microarray data

Martin Schumacher^{a *}, Harald Binder^b, Thomas Gerds^b

^aDepartment of Medical Biometry and Statistics, Institute of Medical Biometry and Medical Informatics, University Medical Center Freiburg, Germany

^b Freiburg Center of Data Analysis and Model Building, University Freiburg, Germany

Addition of a marker to a model

- Typically small improvement in discriminative ability according to c statistic
- c stat blamed for being insensitive

Special Report

Use and Misuse of the Receiver Operating Characteristic Curve in Risk Prediction

Nancy R. Cook, ScD

Letter by Pepe et al Regarding Article, “Use and Misuse of the Receiver Operating Characteristic Curve in Risk Prediction”

To the Editor:

Current statistical approaches for evaluation of risk prediction markers are unsatisfactory. We applaud Cook’s criticisms of the c -index, or area under the receiver operating characteristic curve. This index is based on the notion of pairing subjects, one with poor outcome (eg, cardiovascular event within 10 years) and one without, and determination of whether the risk for the former (ie, the case) is larger than the risk for the latter (ie, the control). This probability of correct ordering of risks is not a relevant measure of clinical value. It should not play a central role in evaluation of risk markers.

Erasmus MC



Editorial JNCI, July 16, 2008 on paper by Gail

Gauging the Performance of SNPs, Biomarkers, and Clinical Factors for Predicting Risk of Breast Cancer

Margaret S. Pepe, Holly E. Janes

Predicting risk of cancer for individuals has long been a goal of medical research. If an individual's risk could be predicted, then prevention and screening modalities could be targeted toward those at meaningfully high risk. This approach is not only more cost efficient than targeting the whole population but also more ethical, at least when interventions are burdensome to the individual. The quest for risk predictors has been revitalized with the emergence of technologies that measure genetic information and other molecular and physiological attributes of the individual. In this issue of the Journal, Gail (1) asks to what extent newly discovered associations between seven single-nucleotide polymorphisms

0.05. Although an extreme example, it illustrates the point. These two criticisms of AUC apply generally, not solely to risk prediction. The AUC really is a poor metric for evaluating markers for disease diagnosis, screening, or prognosis. The third criticism, which is specific to risk prediction, is that the AUC, and indeed the ROC curve itself, hides the values of risk calculated for subjects in the population. Indeed, the risk values are not visible from the ROC curve or the related curves in figure 2 of Gail (1). Moreover, the same ROC curve results if risk values are transformed monotonically, say, multiplied by a factor of 10, yet the clinical implications of these risk values would be very different.

Erasmus MC



Alternatives to ROC analysis

Without harm – benefit: *Stat Med* 2008: 27:157–172;
see S. Greenland commentary

Evaluating the added predictive ability of a new marker: From area under the ROC curve to reclassification and beyond

Michael J. Pencina^{1,*,†}, Ralph B. D’Agostino Sr¹, Ralph B. D’Agostino Jr²
and Ramachandran S. Vasan³

Am J Epidemiol 2008;167:362–368

Practice of Epidemiology

Integrating the Predictiveness of a Marker with Its Performance as a Classifier

Margaret S. Pepe^{1,2}, Ziding Feng¹, Ying Huang², Gary Longton¹, Ross Prentice¹, Ian M. Thompson³, and Yingye Zheng¹

sMC
afung

Alternatives to ROC analysis

With harm – benefit: *Biostatistics* (2008), in press

Estimating the capacity for improvement in risk prediction with a marker

WEN GU, MARGARET SULLIVAN PEPE*

5.4 Risk thresholds and decisions to ascertain Y

Implicit in the above discussion is the existence of threshold values for risk that are used to make decisions, in this case for or against the renal arteriography procedure. Risk thresholds vary with the clinical context and may additionally vary among individuals. How to choose a risk threshold? The classic decision theoretic solution to choosing a risk threshold is fairly simple. Let C_0 (and B_0) denote the cost

ACKNOWLEDGMENTS

We thank Dr A. Cecile J. W. Janssens for allowing us to use the renal artery stenosis data to illustrate our methodology and Dr Patrick Bossuyt for very helpful suggestions.



Contents

1. Developments in performance evaluation
predictions / decisions
2. Evaluation of clinical usefulness
 - A. Binary marker / test
 - B. Additional value of a marker
3. Further developments

Example: Binary markers / tests

- 2 uncorrelated binary markers with equal costs
 - 50% and 10% prevalence, outcome incidence 50%
 - Odds Ratio 4 and 16
 - Evaluate as single test

	Test 1	Test 2
▪ C stat	0.67	0.59
Brier	0.22	0.23
R ²	15%	13%

- Any role for test 2?

Decision threshold and relative costs

Treatment: $\text{Risk} \geq \text{cutoff}$

No treatment: $\text{Risk} < \text{cutoff}$

Event	No event
cTP	cFP
cFN	cTN

cTP and cFP: costs of True and False Positive classifications;
cFN and cTN: costs of False and True Negative classifications respectively.

- Optimal cutoff:
- $\text{Odds}(\text{cutoff}) = (\text{cFP} - \text{cTN}) / (\text{cFN} - \text{cTP})$
= harm / benefit

Simple usefulness measures given 1 cutoff

- Naïve: Unweighted

Sensitivity = $TP / (TP + FN)$; Specificity = $TN / (FP + TN)$

Accuracy: $(TN + TP) / N$; Error rate: $(FN + FP) / N$

Example

- 2 uncorrelated binary markers with equal costs
 - 50% and 10% prevalence, 50% outcome incidence
 - Odds Ratio 4 and 16
 - Evaluate as single test

	Test 1	Test 2
▪ C stat	0.67	0.59
Brier	0.22	0.23
R ²	15%	13%

- Any role for test 2 alone?

Sens	67%	18.8%
Spec	67%	98.7%

Simple usefulness measures given 1 cutoff

- Naïve: Unweighted

Sensitivity = $TP / (TP + FN)$; Specificity = $TN / (FP + TN)$

Accuracy: $(TN + TP) / N$; Error rate: $(FN + FP) / N$

- Weighted variants

Weighted accuracy: $(TP + w TN) / (N_{\text{Event}} + w N_{\text{No event}})$ (Vergouwe 2002)

Net Benefit: $(TP - w FP) / N$,

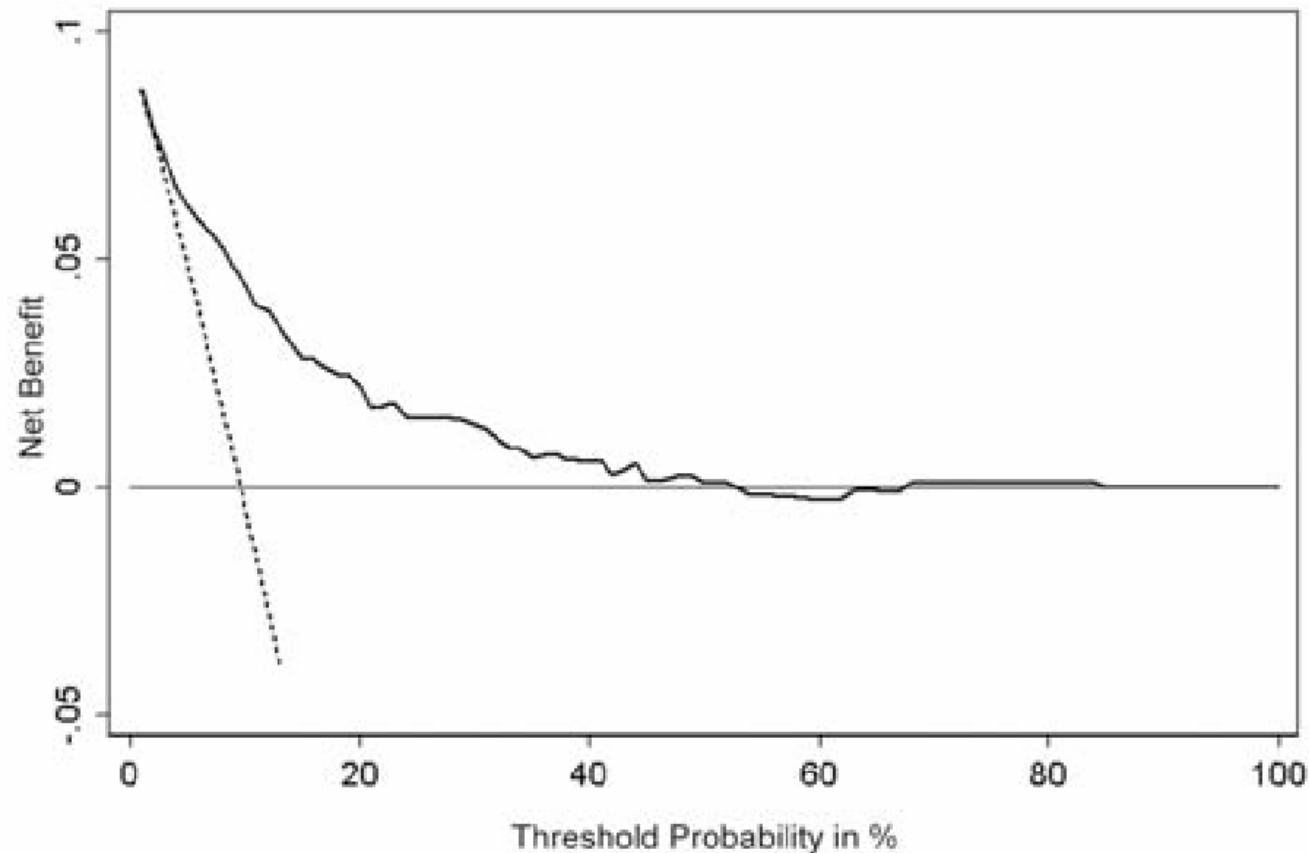
with $w = \text{harm} / \text{benefit}$ (Pierce 1884, Vickers 2006)

From 1 cutoff to consecutive cutoffs

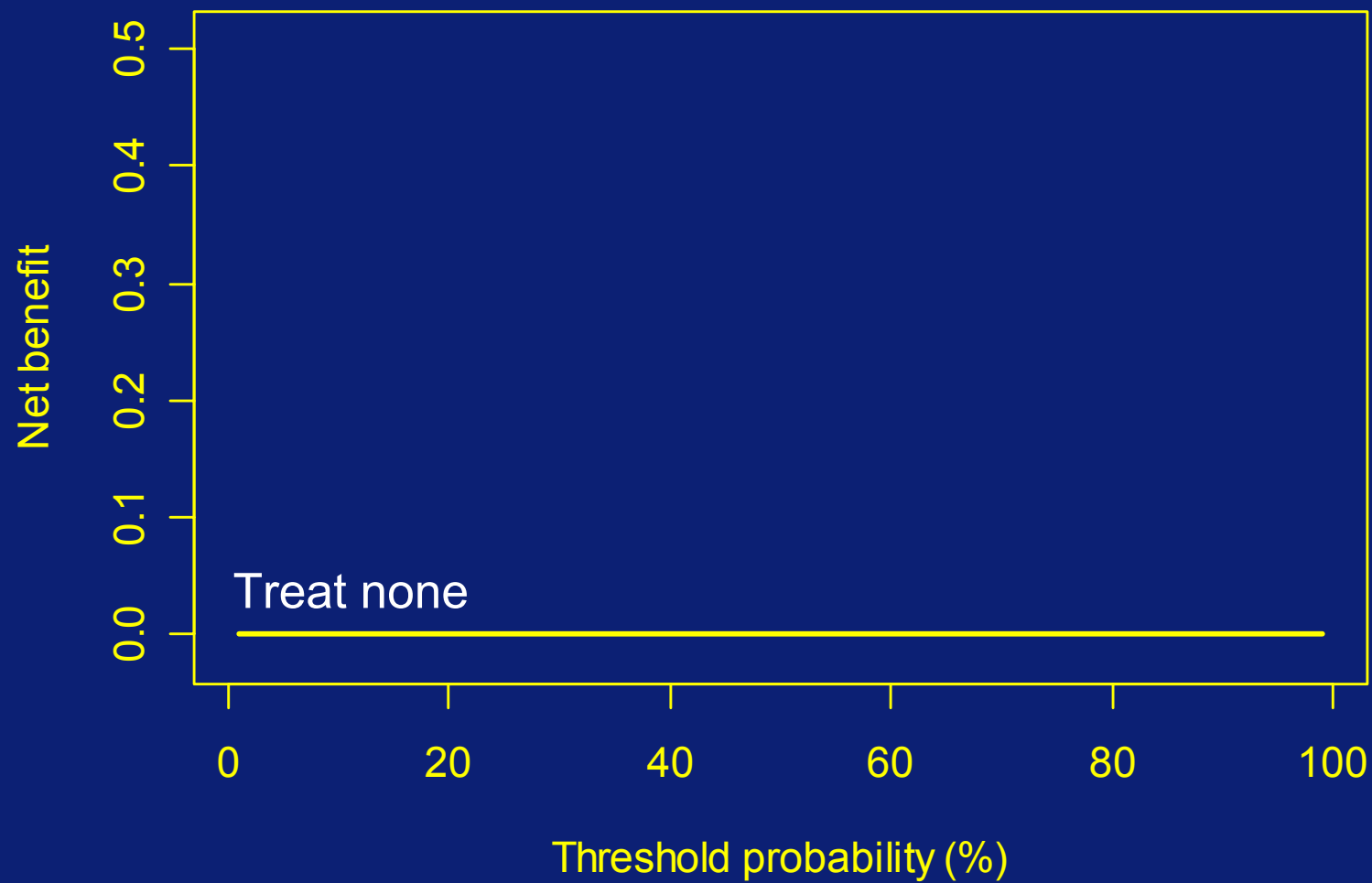
- Sensitivity and specificity → ROC curve
- Net benefit → decision curve

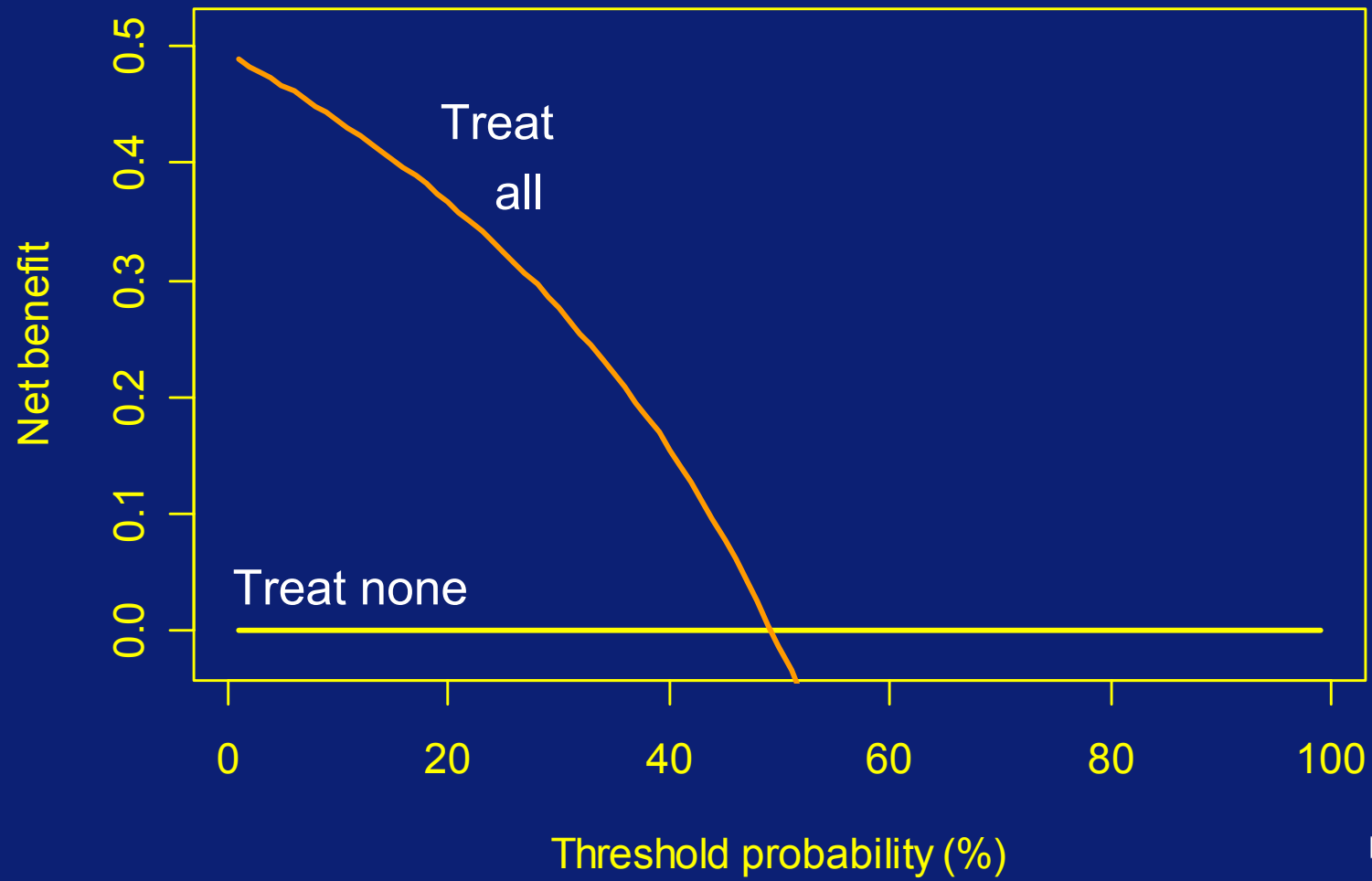
Decision Curve Analysis: A Novel Method for Evaluating Prediction Models

Andrew J. Vickers, PhD, Elena B. Elkin, PhD

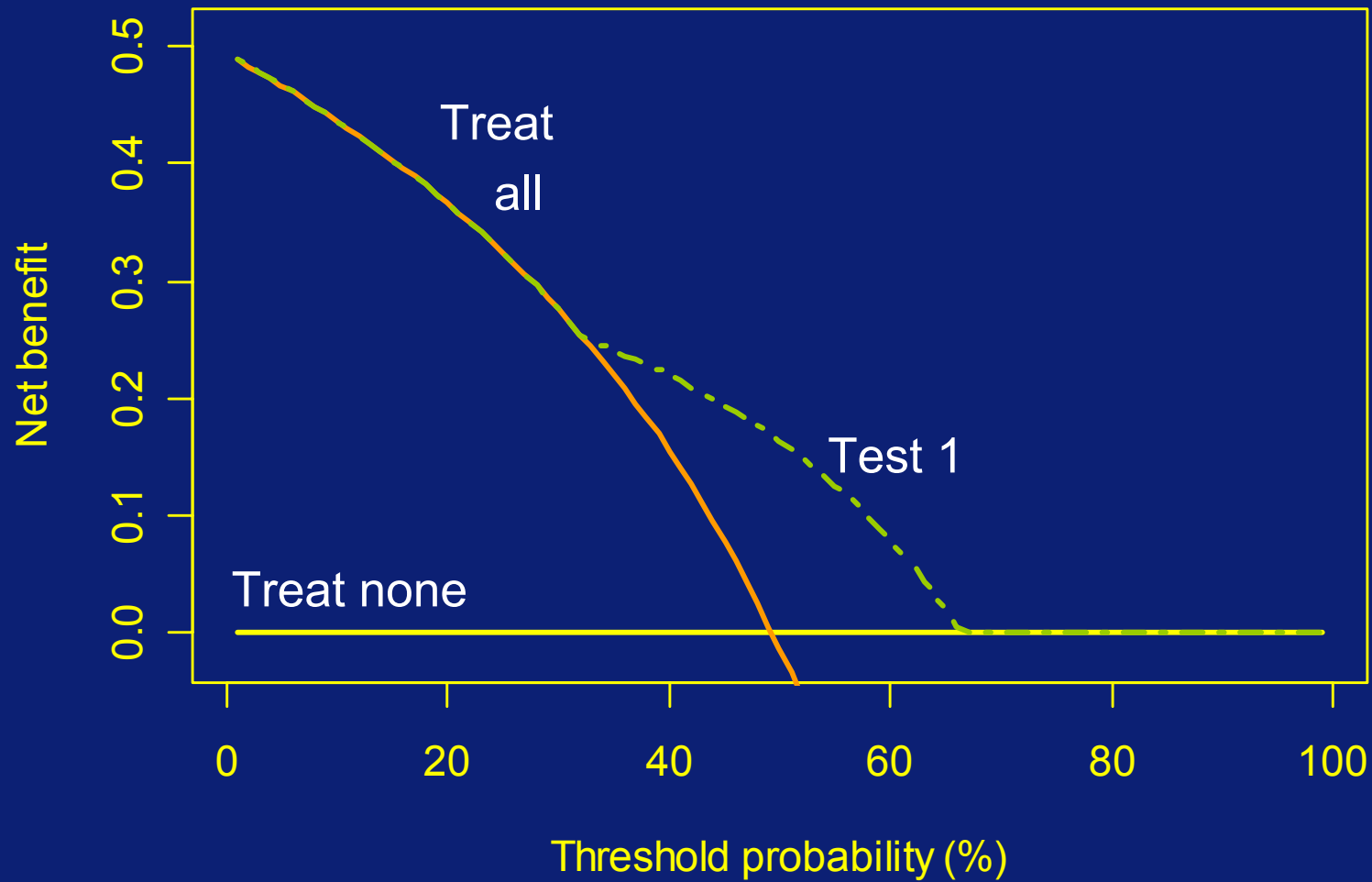


Treat none

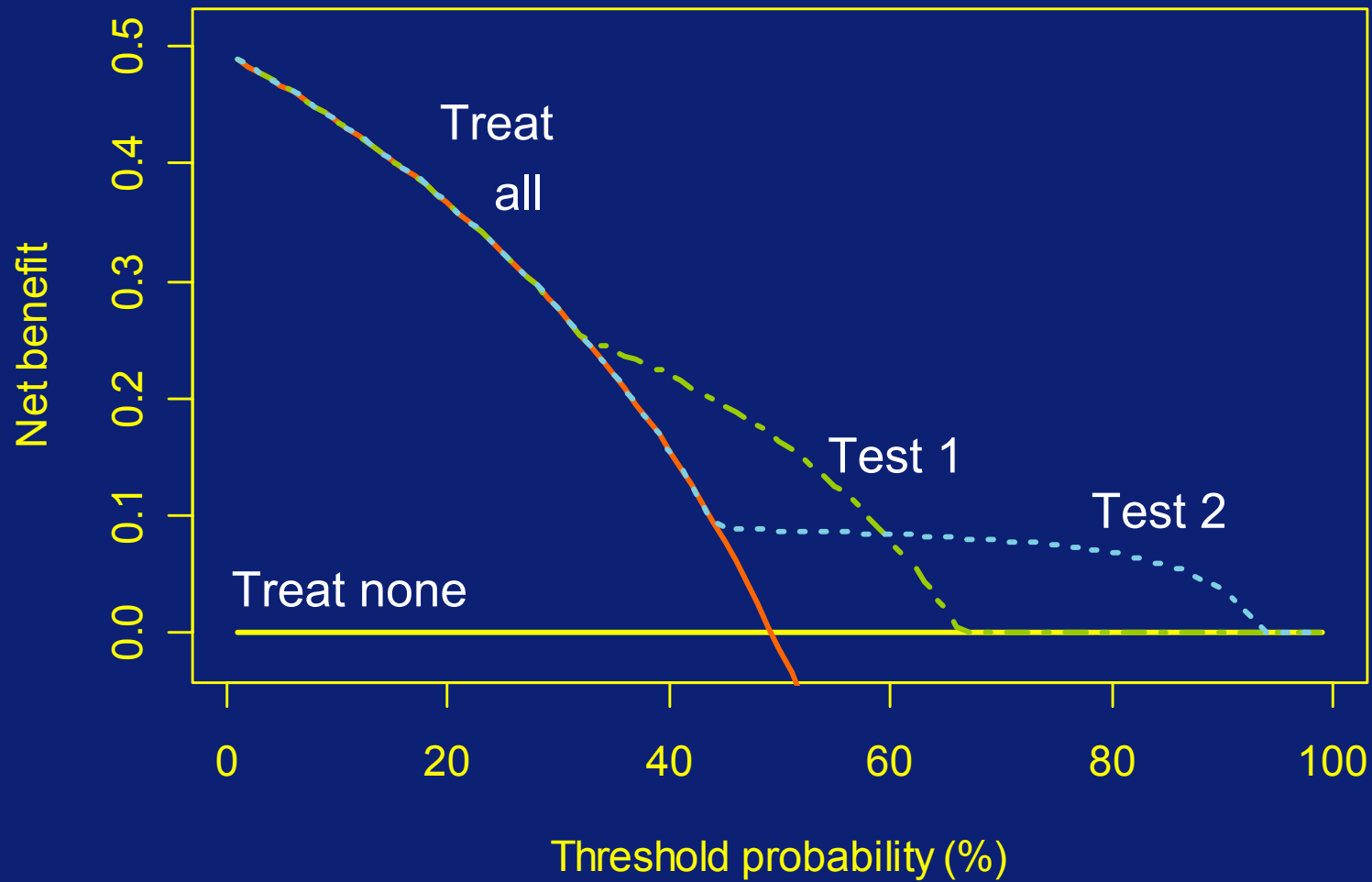




Decision curve for example: Test 1 alone



Decision curve for example: Test 1 and test 2 each alone



Addition of a marker to a model

Additional value of a marker: prostate cancer

- Men with elevated PSA are referred to prostate biopsy
- Only 1 in 4 men with high PSA have prostate cancer
- Could an additional marker help predict biopsy result?
 - Free PSA (a subfraction of PSA)
 - PSA velocity (measure change in PSA)
- Assume decision threshold around 20%

Data set

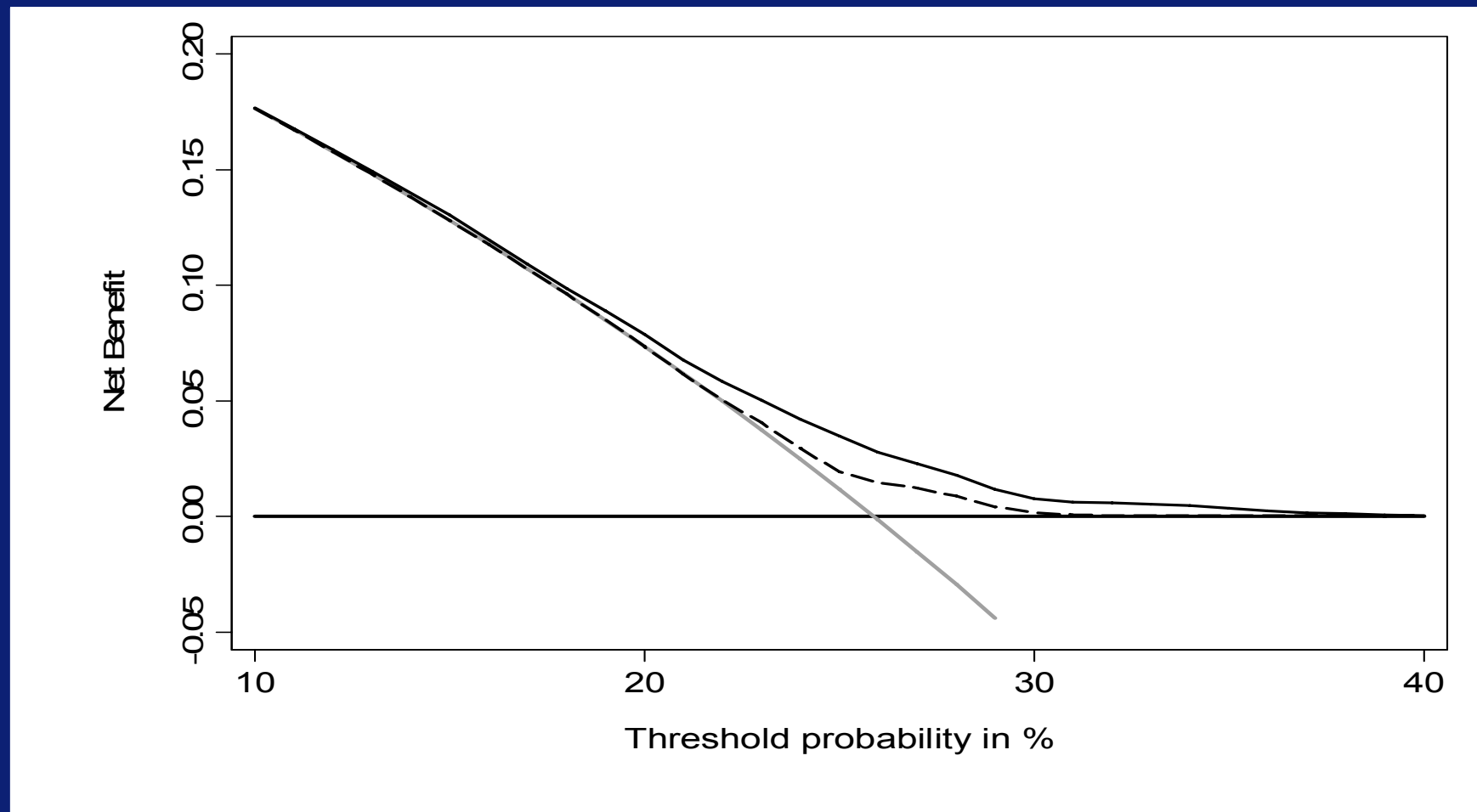
- Data from European Randomized Study of Prostate Cancer screening (ERSPC)
- 2742 previously screened men with:
 - Elevated PSA
 - No previous biopsy
- 710 cancers (26%)

Accuracy metrics

Model	Sens.*	Spec.*	PPV*	NPV*	Brier	AUC	NRI
PSA only	100	0	26	0	.191	.544	
+ PSA velocity	95	10	27	86	.189	.580	.053
+ Free PSA	98	4	26	84	.186	.592	.018
+ Free PSA & PSA velocity	95	8	27	83	.184	.610	.037

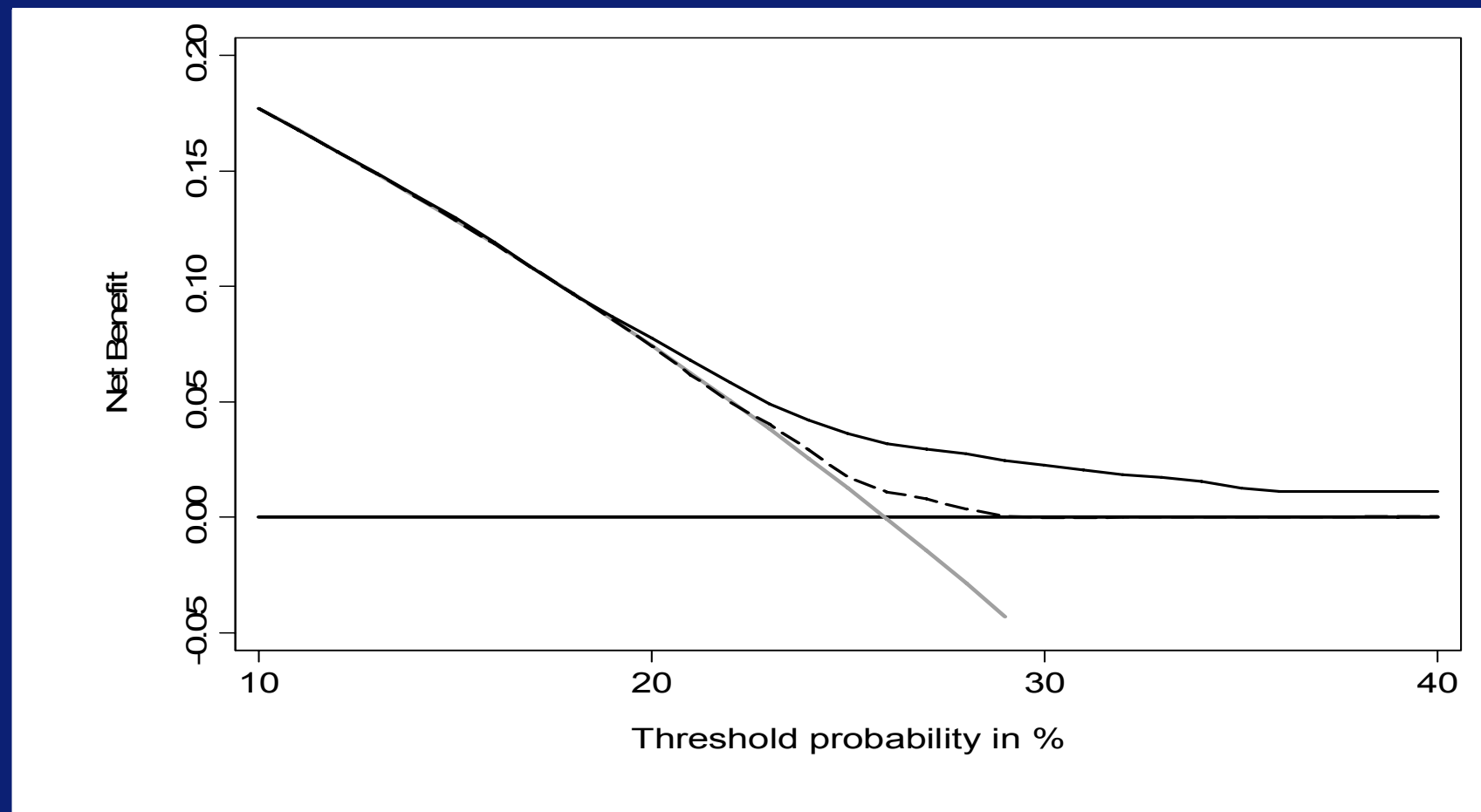
* At Risk threshold of 20%

Add PSA velocity to base model?

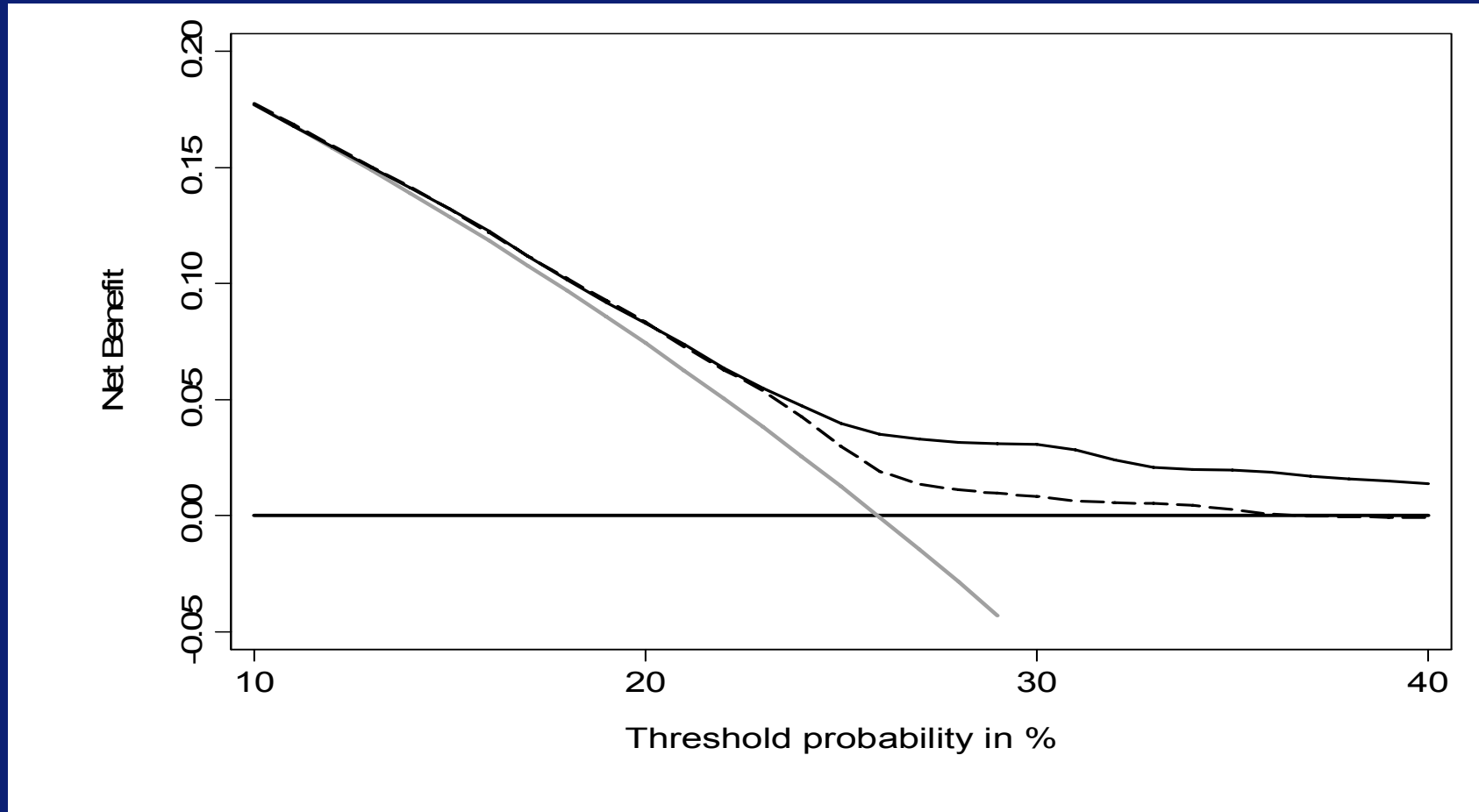


Erasmus

Add free PSA to base model?



Does Free PSA add anything if velocity included?



Accuracy metrics

Model	Sens.*	Spec.*	PPV*	NPV*	Brier	AUC	NRI
PSA only	100	0	26	0	.191	.544	
+ PSA velocity	95	10	27	86	.189	.580	.053
+ Free PSA	98	4	26	84	.186	.592	.018
+ Free PSA & PSA velocity	95	8	27	83	.184	.610	.037

* At Risk threshold of 20%

Which performance measure when?

Application area	Calibration	Discrimination	Clinical usefulness
<i>Public health</i>			
Targeting of preventive interventions			
Predict incident disease	x	X	x
<i>Clinical practice</i>			
Diagnostic work-up			
Test ordering	X	x	X
Starting treatment	X	x	X
Therapeutic decision making			
Surgical decision making	X	x	X
Intensity of treatment	X	x	X
Delaying treatment	X	x	X
<i>Research</i>			
Inclusion in a RCT	X	x	X
Covariate adjustment in a RCT		X	
Confounder adjustment with a propensity score			
Case-mix adjustment			

1. Discrimination: if poor, usefulness unlikely, but NB ≥ 0
2. Calibration: if poor **in new setting**, risk of NB < 0 ;

Prediction model may harm rather than support decision-making

Phases of marker evaluation (Pepe, *Stat Med* 2005;24(24):3687-96)

<i>Phase</i>	<i>Objective</i>	<i>Study design</i>
1 Preclinical exploratory	Promising directions identified	Case-control (convenient samples)
2 Clinical assay and validation	Determine if a clinical assay detects established disease	Case-control (population based)
3 Retrospective longitudinal	Determine if the biomarker detects disease before it becomes clinical. Define a 'screen positive' rule	Nested case-control in a population cohort
4 Prospective screening	Extent and characteristics of disease detected by the test; false referral rate	Cross-sectional population cohort
5 Cancer control	Impact of screening on reducing the burden of disease on the population	Randomized trial

Phases of model development (Reilly *Ann Intern Med* 2006;144(3):201-9)

Level of evidence	Definitions and standards of evaluation	Clinical implications
Level 1		
▪ Derivation of prediction model	▪ Identification of predictors for multivariable model; blinded assessment of outcomes.	▪ Needs validation and further evaluation before using in actual patient care.
Level 2		
▪ Narrow validation of prediction model	▪ Assessment of predictive ability when tested prospectively in 1 setting; blinded assessment of outcomes.	▪ Needs validation in varied settings; may use predictions cautiously in patients similar to sample studied.
Level 3		
▪ Broad validation of prediction model	▪ Assessment of predictive ability in varied settings with wide spectrum of patients and physicians.	▪ Needs impact analysis; may use predictions with confidence in their accuracy.
Level 4		
▪ Narrow impact analysis of prediction model used as decision rule	▪ Prospective demonstration in 1 setting that use of decision rule improves physicians' decisions (quality or cost-effectiveness of patient care).	▪ May use cautiously to inform decisions in settings similar to that studied.
Level 5		
▪ Broad impact analysis of prediction model used as decision rule	▪ Prospective demonstration in varied settings that use of decision rule improves physicians' decisions for wide spectrum of patients.	▪ May use in varied settings with confidence that its use will benefit patient care quality or effectiveness.

Conclusions

- Evaluation of $p(\text{outcome})$ may include overall performance, discrimination and calibration aspects
 - Confusion: overall performance and discrimination measures can be interpreted as evaluation of decision-making
- Evaluation of quality of decision-making requires utility-based loss functions, such as decision-curves

References

- Vickers AJ, Elkin EB: Decision curve analysis: a novel method for evaluating prediction models. *Med Decis Making* 26:565-74, 2006
- Steyerberg EW, Vickers AJ: Decision Curve Analysis: A Discussion. *Med Decis Making* 28; 146, 2008
- Steyerberg EW, et al: Prediction of residual retroperitoneal mass histology after chemotherapy for metastatic nonseminomatous germ cell tumor: analysis of individual patient data from six study groups. *J Clin Oncol* 13:1177-87, 1995
- Vergouwe et al: Predicting retroperitoneal histology in postchemotherapy testicular germ cell cancer: a model update and multicentre validation with more than 1000 patients. *Eur Urol* 51: 424-32, 2007

Read more ...

Books on prediction models

- Cost-effectiveness
 - Costs/page?
 - Costs/formula?
 - Costs/New information
 - Accessibility/Mathematical correctness
- 2 classics + 2 new

Springer Series in Statistics

**Trevor Hastie
Robert Tibshirani
Jerome Friedman**

The Elements of Statistical Learning

**Data Mining, Inference,
and Prediction**



Springer Series in Statistics

Frank E. Harrell, Jr.

Regression Modeling Strategies

**With Applications to
Linear Models,
Logistic Regression,
and Survival Analysis**



WILEY

Multivariable Model-building

A pragmatic approach to regression
analysis based on fractional polynomials
for modelling continuous variables



Patrick Royston and Willi Sauerbrei

WILEY SERIES IN PROBABILITY AND STATISTICS

Erasmus MC

Erasmus

E.Steyerberg@ErasmusMC.nl

Statistics for Biology and Health

Ewout W. Steyerberg

Clinical Prediction Models

A Practical Approach to
Development, Validation, and
Updating



springer.com

20% prepublication discount

Ewout W. Steyerberg, Erasmus University, Rotterdam, The Netherlands

Clinical Prediction Models

A Practical Approach to Development, Validation, and Updating

This book provides a practical approach to the application of modern statistical concepts in clinical prediction problems. The focus is on developing, validating and updating of regression models for diagnosis (presence of disease) and prognosis (outcome of disease). The book includes many examples with R computer code to perform state-of-the-art analyses. The intended readership includes researchers in epidemiology and biostatistics who are involved in prediction modeling, as well as clinicians and policy makers who are interested in learning more about the methodological aspects of prediction models.

Contents: Part I: Prediction models in medicine (Applications • Study design • Statistical models • Overfitting). Part II: Seven steps to develop a valid prediction model (1. Preliminary steps, especially dealing with missing values • 2. Coding of predictors • 3. Model specification • 4. Model estimation • 5. Model performance • 6. Model validation • 7. Model presentation). Part III: Generalizability of prediction models (Patterns of external validity • Updating of prediction models). Part IV: Applications (Case studies for binary and survival outcomes • Data sets for exercises).

2008. Approx. 600 p. Hardcover Statistics for Biology and Health

Publication price: €69.95 **Prepublication price: €55.96**

ISBN: 978-0-387-77243-1

Valid until April 1, 2008. This applies to personal orders only. If you order by credit card, you will not be billed until the book is shipped.

Order Now!

Yes, please send me — copies Statistics for Biology, Health
ISBN: 978-0-387-77243-1 € € € ▶ Prepublication: €55.96

☐ Please bill me
☐ Please charge my credit card: ☐ Eurocard/Access/Mastercard ☐ Visa/Barclaycard/Bank/American ☐ AmericanExpress

Number Valid until

Available from

Springer
Distribution Center GmbH
Haberstr. 7
69126 Heidelberg
Germany

Name	<input type="text"/>
Unit	<input type="text"/>
Institution	<input type="text"/>
Street	<input type="text"/>
City / ZIP-Code	<input type="text"/>
Country	<input type="text"/>
Email	<input type="text"/>
Date	<input type="text"/>
Signature	<input type="text"/>

▶ Call: +49 (0) 6221-345-4301 ▶ Fax: +49 (0) 6221-345-4229
▶ Email: SDC-bookorder@springer.com

If you have any questions, please contact your local Springer office. In Germany, the VAT for books and the VAT for electronic products (digitalization, printing, online information) are not added to the price of the book. Prices are subject to change without notice. All prices are exclusive of carriage charges. Prices and other details are subject to change without notice. All errors and omissions excepted.

Thank you for your attention

Comparison of performance measures

Aspect	Measure	Development *	Validation
Overall performance	R^2	38%	27%
	Brier _{scaled}	28%	20%
Discrimination	C statistic	0.81	0.79
Calibration	Calibration-in-the-large	-	-0.03
	Calibration slope	0.97	0.74
	Test for miscalibration	p=1	p=0.13
Clinical usefulness cutoff 30%	Accuracy	69%	75%
	Net Benefit – resection in all	$0.39 - 0.36 = 0.03$	$0.60 - 0.60 = 0$