

Section III: Evaluating markers

- ▶ Descriptive devices
- ▶ Assessing model calibration
- ▶ Recommended measures of marker performance
 - ▶ Estimation and inference
- ▶ Critique of other marker performance measures
- ▶ Implications for comparing markers or rules
- ▶ Extensions:
 - ▶ Formally incorporating treatment downsides
 - ▶ Evaluating a prognostic marker

Descriptive devices

- ▶ Risk curves
- ▶ Treatment effect curves

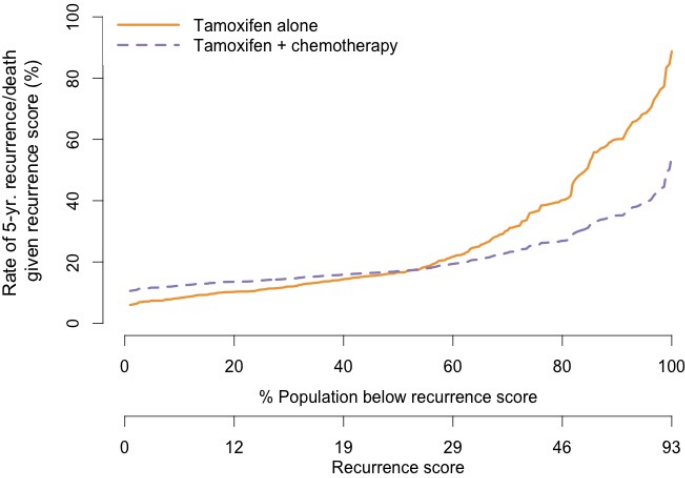
Terminology suggests the outcome is binary, but these devices also apply to categorical and continuous outcomes.

Risk curves

For a single marker, *risk curves* plot the expected outcome as a function of the marker, for each treatment.

We recommend aligning the curves for the two treatment groups with respect to marker percentile $F(X)$, rather than marker value X . I.e., plot $E(D|A, X)$ vs. $F(X)$ for $A = 0, 1$.

Example: Oncotype DX marker in the breast cancer trial



Treatment effect curves

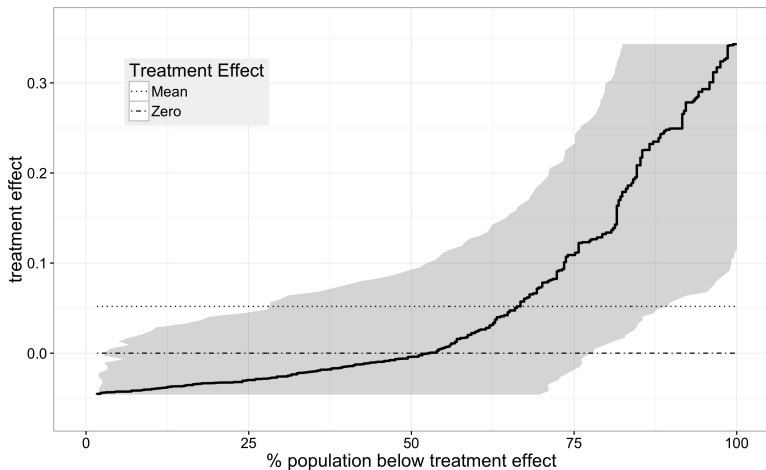
Show the distribution of the marker-specific treatment effect, $\Delta(X) = E(D|A = 0, X) - E(D|A = 1, X)$.

Different scales are possible:

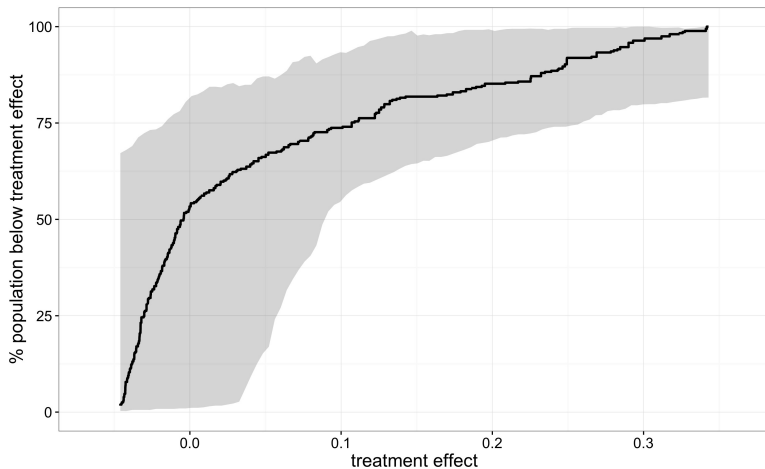
- ▶ Reverse-CDF, i.e. $\Delta(X) = \delta$ vs. $F_{\Delta}(\delta)$. Also called a *predictiveness curve* (Huang et al. 2007).
- ▶ Traditional CDF, i.e. $F_{\Delta}(\delta)$ vs. $\Delta(X) = \delta$.
- ▶ Density or histogram of $\Delta(X)$.

Unlike the risk curve plot, this device applies to multivariate X .

Treatment effect curve for the Oncotype DX marker: Reverse CDF



Treatment effect curve for the Oncotype DX marker: Traditional CDF



Checking model calibration

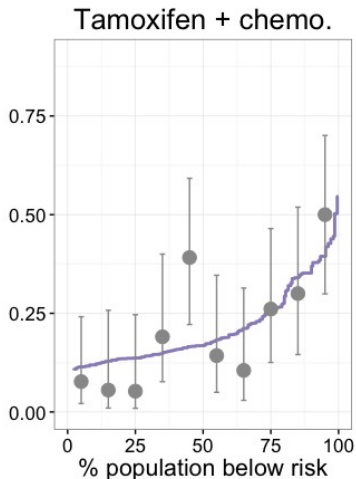
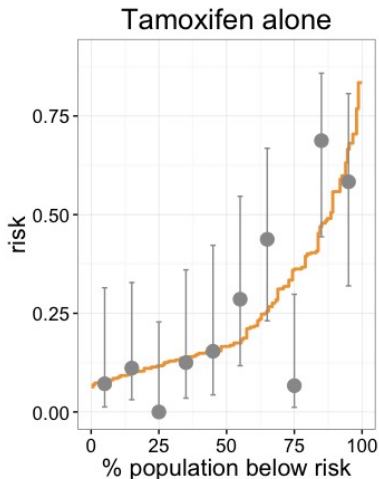
Estimating these curves requires modeling $E(D|A, X)$.

Good calibration of the $E(D|A, X)$ model is essential for validity of the risk and treatment effect curves.

Two approaches to assessing calibration:

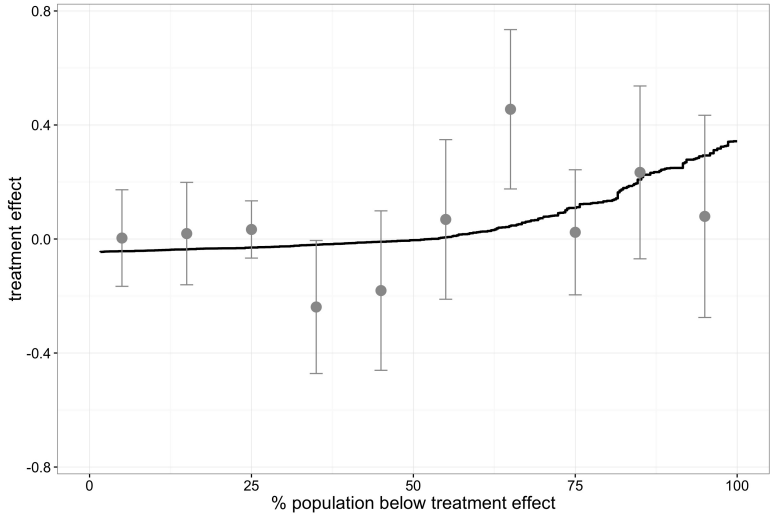
- ▶ Overlay observed risks and treatment effects on the plots
- ▶ Formally compare observed vs. predicted values using Hosmer-Lemeshow goodness of fit tests

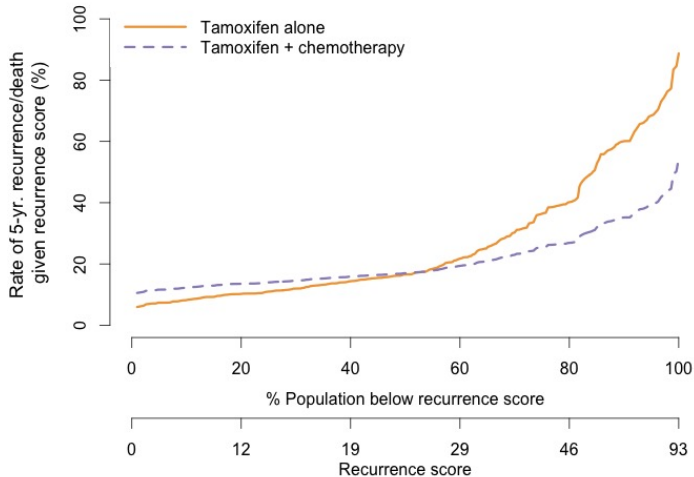
Example: Oncotype DX risk curve calibration



No significant difference between observed and predicted risks in either treatment group ($p = 0.078$ and 0.096 , Hosmer-Lemeshow test).

Example: Oncotype DX treatment effect curve calibration





Is the marker good enough to incorporate into clinical practice?

Performance Measures

Context

Goal is to evaluate the performance of marker-based treatment rule $d(X)$.

- ▶ a rule estimated from the data, $d_n(X)$, or a pre-specified rule

Focus on the setting where $A = 0$ is the default treatment choice absent X .

- ▶ X is used to identify a subgroup likely to benefit from treatment

The opposite scenario where $A = 1$ is the default and X is used to identify a subgroup not likely to benefit from treatment is handled analogously.

Evaluating a marker-based treatment rule

Suppose that $A = 0$ is the default choice absent X .

The *clinical impact* of rule $d(X)$ is

$$\begin{aligned}\mathcal{I}(d) &= E(D \mid A = 0) - E(D \mid \text{treat using rule } d) \\ &= E(D \mid A = 0) - E(D(d)) \\ &= [E(D \mid A = 0, d(X) = 0) \cdot P(d(X) = 0) \\ &\quad + E(D \mid A = 0, d(X) = 1) \cdot P(d(X) = 1)] \\ &\quad - [E(D \mid A = 0, d(X) = 0) \cdot P(d(X) = 0) \\ &\quad + E(D \mid A = 1, d(X) = 1) \cdot P(d(X) = 1)] \\ &= E(\Delta(X) \mid d(X) = 1) \cdot P(d(X) = 1) \\ &\equiv \beta(d) \cdot \tau(d)\end{aligned}$$

The two constituents of $\mathcal{I}(d)$,

- ▶ $\tau(d)$ = the proportion of subjects impacted by X measurement, who are recommended treatment
- ▶ $\beta(d)$ = average treatment efficacy in this subgroup

are important measures in their own right.

In practice we recommend reporting the triplet $(\mathcal{I}(d), \tau(d), \beta(d))$, along with the expected outcomes under “treat all” and “treat none” policies, $\rho_0 = E(D|A = 0)$ and $\rho_1 = E(D|A = 1)$.

If $A = 1$ is the default choice absent X ,

$$\begin{aligned}\mathcal{I}(d) &= E(D \mid A = 1) - E(D \mid \text{treat using rule } d) \\ &= E(-\Delta(X) \mid d(X) = 0) \cdot P(d(X) = 0)\end{aligned}$$

and the two constituents are

- ▶ $\tau(d) = P(d(X) = 0)$ = the proportion impacted by X measurement, who are recommended no treatment
- ▶ $\beta(d) = E(-\Delta(X) \mid d(X) = 0)$ = average benefit of no treatment in this subgroup

Empirical estimation

Estimate the performance of rule $d(X)$ empirically using

$$\hat{\tau}^e(d) = \mathbb{P}(d(X) = 1)$$

$$\hat{\beta}^e(d) = \mathbb{E}(D \mid A = 0, d(X) = 1) - \mathbb{E}(D \mid A = 1, d(X) = 1)$$

$$\hat{\mathcal{I}}^e(d) = \hat{\beta}^e(d) \cdot \hat{\tau}^e(d)$$

where \mathbb{P} is the empirical probability and \mathbb{E} is the empirical mean

- ▶ Equivalent to estimating $E(D(d))$ using the IPW estimator defined in Section II

Model-based estimation

When $d(X)$ is derived using a model for $E(D|A, X)$, the resultant model for $\Delta(X)$ can be used to estimate performance in a *model-based* fashion:

$$\begin{aligned}\hat{\beta}^m(d) &= \mathbb{E}(\hat{\Delta}(X) \mid d(X) = 1) \\ \hat{\mathcal{I}}^m(d) &= \mathbb{E}(\hat{\Delta}(X) I(d(X) = 1)) \\ &= \hat{\beta}^m(d) \cdot \hat{\tau}^e(d)\end{aligned}$$

Model-based estimators are more efficient. However they are biased if the $E(D|A, X)$ model is mis-specified.

Inference

When evaluating performance of a pre-specified rule $d(X)$, all estimates of performance are asymptotically normal. Quantile bootstrap confidence intervals work well.

Similarly, when training data are used to derive $d_n(X)$ and independent test data are used to estimate performance, estimators are asymptotically normal and the bootstrap can be used for inference.

One exception to the above is when $P(\Delta(X) = 0) > 0$, i.e. there exist subjects with $\Delta(X)$ identically 0. Performance estimates may not be asymptotically normal and the bootstrap may not perform well.

Inference, continued

However, when the same data are used to derive $d_n(X)$ and to estimate performance

- ▶ Estimates are biased (overoptimistic)
- ▶ Estimators are not asymptotically normal. Bootstrap-based confidence intervals may not have good coverage.
- ▶ Performance of normal-theory/bootstrap inferential methods is expected to be worse for settings with: small n , high-dimensional X , heavy marker/model selection
- ▶ There are partial solutions (next slide).
- ▶ This is an active research area.

Partial solutions to drawing inference absent test data

Cross-validation (CV)

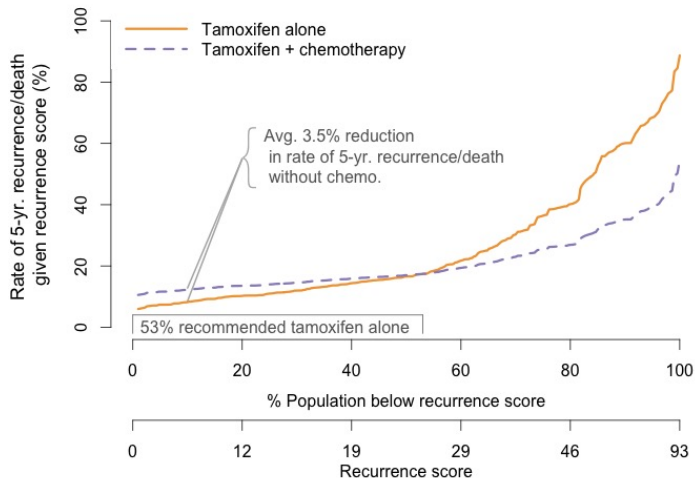
- ▶ Sample B training/test data splits. For each, obtain $d_n^b(X)$ using training data and estimate performance using test data. Average performance estimates. Shift naive performance estimates and confidence intervals down by the estimated bias – the difference between naive and CV performance estimates.

Bootstrap bias correction: the “refined bootstrap” (Efron and Tibshirani 1994)

- ▶ Sample B bootstrap datasets. For each, obtain $d_n^b(X)$ and calculate the difference in estimated performance of this rule using the bootstrap vs. original data. The average of these differences estimates the bias. Shift naive performance estimates and confidence intervals down by the estimated bias.

There are variations on each of these approaches.

Example: Oncotype DX marker performance



Risk curves estimated using logistic regression. Bootstrap-bias-corrected empirical performance estimates.

Absent X , chemotherapy is the default.

Given X ,

- ▶ $\hat{\tau} = 53.0\%$ avoid chemo, and associated toxicity and cost (0.2 to 80.1)
- ▶ $\hat{\beta} = 3.5\%$ lower risk of 5-yr. recurrence/death in subset avoiding chemo. (-12.9 to 10.8)
- ▶ Estimated clinical impact is $\hat{\mathcal{I}} = 1.5\%$ lower 5-yr. recurrence/death rate (-3.6 5.7)
 - ▶ 21% event rate under default “chemo. for all” policy is reduced to 19.5% with use of X .
 - ▶ 25% event rate under “chemo for none”

Said another way,

- ▶ The overall efficacy of chemo. is a 3.9% absolute reduction in the 5-yr. recurrence/death rate.
- ▶ The efficacy of X-based chemo. is a $3.9 + 1.5 = 5.4\%$ reduction in the 5-yr. recurrence/death rate.

Other Performance Measures

A marker-by-treatment interaction is insufficient

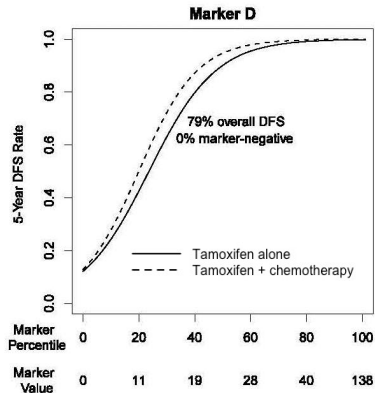
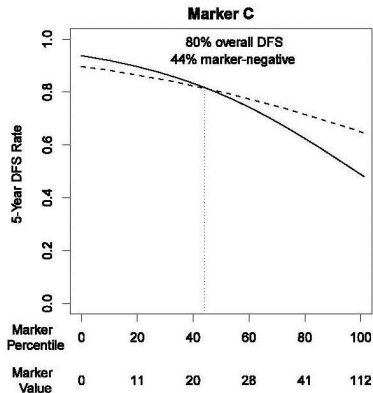
Testing for a marker-by-treatment interaction is a useful first step.

- ▶ An interaction is necessary, but not sufficient, for the marker to have value

However, the interaction coefficient does not *quantify* marker performance.

- ▶ Interpretation depends on the scale of the $E(D|A, X)$ model, the other variables in the model, and the scale of the marker
- ▶ Easy to construct examples of markers with the same interaction coefficient, but different clinical impact

Example



Two markers with the same marker-by-treatment interaction, but very different performance.

Janes et al. (*Ann. Int. Med.* 2011)

What about biomarker accuracy?

Sensitivity, specificity, PPV, and NPV are classic performance measures for diagnostic, screening, and prognostic markers.

FDA and IOM biomarker development guidance documents advocate reporting accuracy measures

- ▶ without properly distinguishing between approaches for diagnostic and prognostic and predictive/treatment selection markers

Accuracy measures have been proposed for treatment selection markers, for the setting of a binary outcome D

Huang et al. (*Biometrics* 2012), Zhang et al. (*Ann Appl Stat* 2014), Sitlani and Heagerty (*Stat Med* 2014), Simon (*JNCI* 2015)

Accuracy measures rely on potential outcomes

$D(0)$ = potential outcome without treatment

$D(1)$ = potential outcome with treatment

Trt. benefit $\equiv D(0) = 1, D(1) = 0$

No trt. benefit $\equiv D(0) = D(1)$ or $D(0) = 0, D(1) = 1$

The accuracy of rule $d(X)$ is then measured by:

Sensitivity = $P(d(X) = 1 \mid \textit{Trt. benefit})$

Specificity = $P(d(X) = 0 \mid \textit{No trt. benefit})$

PPV = $P(\textit{Trt. benefit} \mid d(X) = 1)$

NPV = $P(\textit{No trt. benefit} \mid d(X) = 0)$

Fundamental problem

Almost never can both potential outcomes be observed and so we do not know whether a subject benefits from treatment.

Therefore, in general the accuracy measures are not estimable from data– even RCT data.

Illustration: Two binary markers in an RCT (n = 2000)

Unobservable data: Marker-positivity by potential outcome

		Benefit from trt. (n = 400)	Bad outcome regardless of trt. (n = 600)	Good outcome regardless of trt. (n = 600)	Harmed by trt. (n = 400)
Marker 1	Negative	200	250	400	250
	Positive	200	350	200	150
Marker 2	Negative	100	350	500	150
	Positive	300	250	100	250

Observable data: Marker-positivity by observed outcome

		Treatment arm		No trt. (n = 1000)		Trt. (n = 1000)	
		Outcome		Good	Bad	Good	Bad
Marker 1	Negative	325	225	300	250		
	Positive	175	275	200	250		
Marker 2	Negative	325	225	300	250		
	Positive	175	275	200	250		

The biomarkers have very different accuracy, but the same observed data:

Marker 1 Sensitivity = 50% Specificity = 56%
 PPV = 22% NPV = 82%
 Prop. marker-positive = 56%

Marker 2 Sensitivity = 75% Specificity = 63%
 PPV = 33% NPV = 91%
 Prop. marker-positive = 56%

“Pragmatic” accuracy measures have been proposed which assume $D(0) \perp D(1)$ given X .

- ▶ This assumption is unlikely to hold in any clinical context
- ▶ This example illustrates the fallacy of these pragmatic measures

Pragmatic Accuracy (Both Markers)	Sensitivity _i = 61% PPV _i = 27%	Specificity _i = 46% NPV _i = 78%
--------------------------------------	--	--

Marker 1 Truth	Sensitivity = 50% PPV = 22%	Specificity = 56% NPV = 82%
----------------	--------------------------------	--------------------------------

Marker 2 Truth	Sensitivity = 75% PPV = 33%	Specificity = 63% NPV = 91%
----------------	--------------------------------	--------------------------------

Our recommendation

In general, accuracy estimates depend on unverifiable assumptions about the joint distribution of potential outcomes.

We recommend instead focusing on identifiable marker performance measures: \mathcal{I} , τ , β .

These measures do not depend on the joint distribution of potential outcomes.

Other performance measures

Recall $\rho_0 = E(D|A = 0)$ and $\rho_1 = E(D|A = 1)$. Note that $E(\Delta(X)) = \rho_0 - \rho_1$.

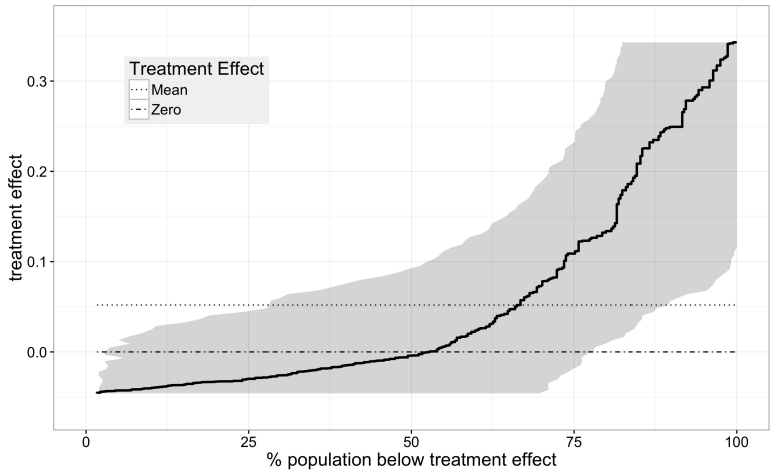
Variance in treatment effect,

$$V_{\Delta} \equiv \int (\Delta(X) - (\rho_0 - \rho_1))^2 \partial F_{\Delta}$$

Total gain,

$$\text{TG} \equiv \int |\Delta(X) - (\rho_0 - \rho_1)| \partial F_{\Delta}$$

- ▶ Two “global” performance measures– do not require specifying a treatment rule
- ▶ They lack a clinically relevant interpretation



Any one performance measure is insufficient

We advocate reporting the triplet

$$\tau(d) = P(d(X) = 1)$$

$$\beta(d) = E(\Delta(X) \mid d(X) = 1)$$

$$\mathcal{I}(d) = \beta(d) \cdot \tau(d)$$

No single measure says it all.

- ▶ E.g., a large β may not be compelling if τ is small
- ▶ E.g., if treatment has downsides not captured in D , \mathcal{I} is insufficient and we need τ to capture treatment “cost”

Implications for comparing markers or treatment rules

Estimate contrasts in the above performance measures.

Again, our recommendation is to contrast

$$\tau(d) = P(d(X) = 1)$$

$$\beta(d) = E(\Delta(X) \mid d(X) = 1)$$

$$\mathcal{I}(d) = \beta(d) \cdot \tau(d)$$

Performance measures can be compared using Wald-type hypothesis tests.

Example: Two simulated markers in the breast cancer context

	Marker X_1 Estimate (95% CI)	Marker X_2 Estimate (95% CI)	X_1 vs. X_2 Estimated Diff. (95% CI)	P-value
$\hat{\tau}^e$	0.461 (0.000,0.700)	0.377 (0.304,0.470)	0.084 (-0.358,0.236)	0.768
$\hat{\beta}^e$	0.029 (-0.106,0.082)	0.238 (0.170,0.309)	-0.209 (-0.342,-0.129)	< 0.002
$\hat{\beta}^m$	0.023 (0.000,0.057)	0.262 (0.209,0.310)	-0.239 (-0.294,-0.178)	< 0.002
\hat{I}^e	0.013 (-0.010,0.044)	0.090 (0.060,0.122)	-0.076 (-0.111,-0.042)	< 0.002
\hat{I}^m	0.010 (0.000,0.037)	0.099 (0.071,0.129)	-0.088 (-0.115,-0.061)	< 0.002

Formally incorporating treatment downsides

Incorporating treatment downsides into the treatment rule

The rule

$$d(X) = I(E(D|A = 0, X) - E(D|A = 1, X) > 0)$$

is optimal if the goal is to minimize $E(D(d))$.

If, however, there are additional downsides of treatment not captured in D , it is compelling to consider rules of the form

$$d^\delta(X) = I(E(D|A = 0, X) - E(D|A = 1, X) > \delta),$$

for $\delta > 0$.

- ▶ Such rules are optimal for maximizing the *net benefit* of marker-based treatment (Vickers et al. 2007; Janes et al. 2014)

Choice of treatment effect threshold, δ

Decision theory suggests that δ should correspond to the cost/dis-utility of treatment A relative to the cost/dis-utility of one unit of the outcome D

- ▶ E.g. if treatment-associated toxicity is 1/10 the cost of the binary clinical outcome, the optimal $\delta = 0.10$

“Cost” is used broadly here; units may be dollars or quality-adjusted life years (QALYs) or probabilities of downstream events.

Choice of treatment effect threshold, continued

This result can be used in the other direction: given δ , the relative importance of D and A (the cost ratio) is quantified.

- ▶ E.g., if $\Delta(X) > 0.02$ justifies a treatment recommendation, this implies that the cost of the binary clinical outcome is 50 times the cost of treatment.

Evaluating performance

The performance of treatment rule $d^\delta(X)$ can be evaluated using the aforementioned metrics:

$$\tau(d^\delta) = P(d^\delta(X) = 1)$$

$$\beta(d^\delta) = E(\Delta(X) \mid d^\delta(X) = 1)$$

$$\begin{aligned}\mathcal{I}(d^\delta) &= E(D \mid A = 0) - E(D \mid \text{treat using rule } d^\delta) \\ &= \beta(d^\delta) \cdot \tau(d^\delta)\end{aligned}$$

Net benefit

The *net benefit* measure captures the net impact of using the marker to select treatment, including its impact on outcomes and on treatment.

Let C_D be the cost/dis-utility of one unit of D and C_A be the cost/dis-utility of treatment.

$$\begin{aligned}\text{NB}(d^\delta) &\equiv \text{Exp. cost under } A = 0 - \text{Exp. cost using } d^\delta \\ &= [E(D \mid A = 0) - E(D \mid \text{using rule } d^\delta)]C_D \\ &\quad - P(d^\delta(X) = 1)C_A \\ &= \mathcal{I}(d^\delta)C_D - \tau(d^\delta)C_A\end{aligned}$$

- ▶ This could be reduced by the dis-utility of measuring the marker in everyone

If the optimal rule for maximizing net benefit is used ($\delta = \frac{C_A}{C_D}$), the net benefit in C_D units is

$$\text{NB}(d^\delta) = \mathcal{I}(d^\delta) - \tau(d^\delta)\delta$$

- ▶ Appealing that this NB formulation depends only on δ , and not on C_D or C_A

Thus, $\text{NB}(d^\delta)$ can be interpreted as the *discounted reduction in the expected outcome* under marker-based treatment.

Note that if $\delta = 0$, $\text{NB}(d^\delta) = \mathcal{I}(d^\delta)$.

If $A = 1$ is the default absent X ,

$$\begin{aligned}\text{NB}(d^\delta) &\equiv \text{Exp. cost under } A = 1 - \text{Exp. cost using } d^\delta \\ &= [E(D | A = 1) - E(D | \text{using rule } d^\delta)]C_D \\ &\quad + [1 - P(d^\delta(X) = 1)]C_A \\ &= \mathcal{I}(d^\delta)C_D + \tau(d^\delta)C_A \\ &= \mathcal{I}(d^\delta) + \tau(d^\delta)\delta,\end{aligned}$$

where the last line holds if $\delta = C_A/C_D$ and NB is in C_D units.

Thus, $\text{NB}(d^\delta)$ can be interpreted as the *augmented reduction in the expected outcome* under marker-based treatment.

Example: Oncotype DX marker performance

Suppose it is determined that $\Delta(X) > 0.01$ is large enough to warrant a chemotherapy recommendation; women with $\Delta(X) < 0.01$ should be recommended no chemo.

We estimate that using this rule to recommend no chemo. would:

- ▶ Allow 56.3% of women to avoid chemo.
- ▶ Reduce the 5-yr. recurrence/death rate in this subgroup by 1.7%
- ▶ Reduce the population 5-yr. recurrence/death rate by 0.5%
- ▶ Yield a NB of 0.011. Thus, 1.1% is the augmented reduction in the 5-yr. recurrence/death rate.

Evaluating a prognostic marker

Definition

A *prognostic* marker/model predicts outcomes under standard of care.

Prognostic markers are often used to guide treatment.

- ▶ E.g. Gail model for predicting breast cancer risk, used to guide use of tamoxifen
- ▶ E.g. Framingham model for predicting CVD risk, used to guide use of statins
- ▶ E.g. Partin tables used to guide treatment of prostate cancer

The (implicit) logic is that subjects at higher risk have a greater absolute benefit from treatment if treatment has a constant relative risk.

Common approach to prognostic marker evaluation

Use the marker to estimate a “risk score”, $E(D | A = 0, X)$, and evaluate the ROC curve (AUC) or (Sens, Spec) of the risk score using a chosen “high risk” threshold.

Or, for contrasting two models $E(D | A = 0, X)$ vs. $E(D | A = 0, Y)$, evaluate the difference in AUCs, difference in (Sens, Spec), or NRI.

These measures do not capture the value of the marker for *guiding treatment*. To this end, we should evaluate the marker in an RCT using the aforementioned performance metrics.

Example: HIV prevention trial

RCT of PrEP vs. placebo for prevention of HIV infection in MSM.

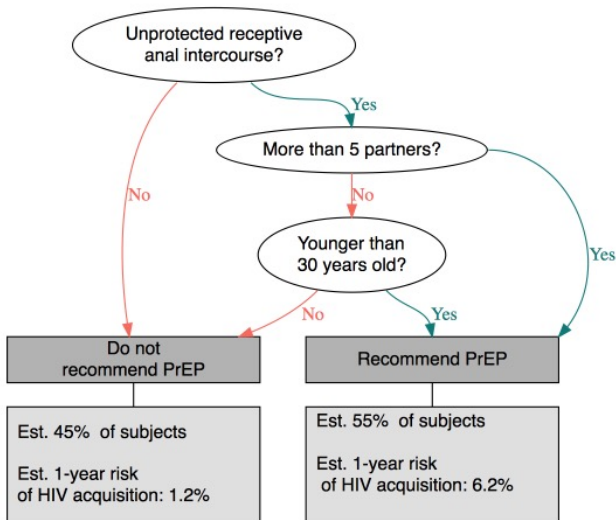
An HIV risk prediction model was developed using data from the placebo arm.

- ▶ Cox proportional hazards logistic regression (Ruczinski et al. 2003)

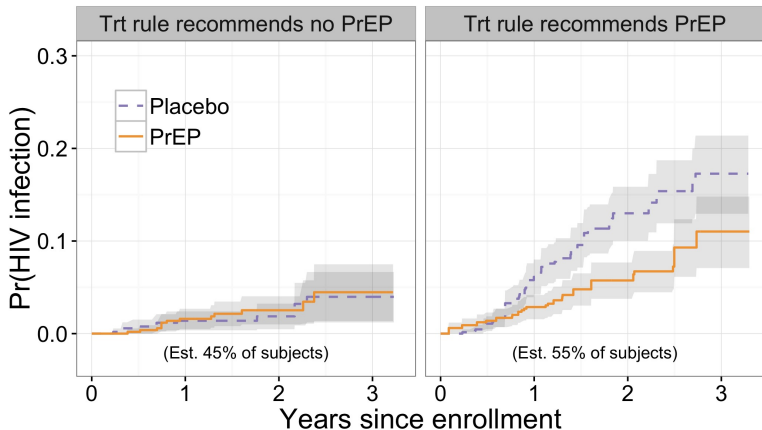
What is the performance of the risk model for guiding treatment?

- ▶ Current WHO guideline recommends PrEP for subjects estimated to be at or above 3% 1-yr. risk without PrEP

Risk-based treatment recommendation



Based on Cox logic regression model fit using placebo-arm data.



Without PrEP, est. 1-yr HIV incidence is 4.0% (2.9 - 5.2%)

PrEP for all yields est. incidence 2.3% (1.4 - 3.2%)

PrEP for high risk subjects yields est. incidence 2.4% (1.8 - 2.9%),
and requires treating only 55.2% (23.4 - 79.2%)

Summary

- ▶ Descriptive devices are useful for visualizing data
- ▶ Clinical impact and its constituents are recommended performance measures
- ▶ Contrasts in these measures are recommended for comparing markers or rules
- ▶ Extensions allow treatment downsides to be incorporated into the treatment rule and its evaluation
- ▶ Prognostic markers used to select treatment should be similarly evaluated