

# Individualized treatment effect was predicted best by modeling baseline risk in interaction with treatment assignment

Alexandros Rekkas<sup>a</sup>, Peter R. Rijnbeek<sup>a</sup>, David M. Kent<sup>b</sup>, Ewout W. Steyerberg<sup>c</sup>, David van Klaveren<sup>d</sup>

<sup>a</sup>*Department of Medical Informatics, Erasmus Medical Center, Rotterdam, The Netherlands*

<sup>b</sup>*Predictive Analytics and Comparative Effectiveness Center, Institute for Clinical Research and Health Policy Studies, Tufts Medical Center, Boston, Massachusetts, USA*

<sup>c</sup>*Department of Biomedical Data Sciences, Leiden University Medical Center, Leiden, The Netherlands*

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

---

## Abstract

**Objective:** To compare different risk-based methods for optimal prediction of individualized treatment effects.

**Study Design and Setting:** We simulated RCT data using diverse assumptions for the average treatment effect, a baseline prognostic index of risk (PI), the shape of its interaction with treatment (none, linear, quadratic or non-monotonic), and the magnitude of treatment-related harms (none or constant independent of the PI). We predicted absolute benefit using: models with a constant relative treatment effect; stratification in quarters of the PI; models including a linear interaction of treatment with the PI; models including an interaction of treatment with a restricted cubic spline (RCS) transformation of the PI; an adaptive approach using Akaike's Information Criterion. We evaluated predictive performance using root mean squared error and measures of discrimination and calibration for benefit. **Results:** The linear-interaction model and the RCS-interaction displayed robust performance across many simulation scenarios. The RCS-model was optimal when quadratic or non-monotonic deviations from a constant treatment effect were stronger, and when sample size was larger. The adaptive approach required larger sample sizes. Illustrations in the GUSTO-I trial confirmed these findings. **Conclusion:** An interaction between baseline risk and treatment assignment should be considered to improve treatment effect predictions.

**Keywords:** treatment effect heterogeneity absolute benefit prediction models

---

## <sup>1</sup> What is new?

- <sup>2</sup> ▪ Models including a linear interaction of the prognostic index with treatment were a good option for predicting
- <sup>3</sup>     absolute benefit with smaller sample sizes or moderately performing prediction models.
- <sup>4</sup> ▪ Models including an interaction of treatment with a restricted cubic spline transformation using 3 knots
- <sup>5</sup>     performed better in the presence of true non-monotonic deviations from a constant relative treatment effect
- <sup>6</sup>     or with treatment-related harms.
- <sup>7</sup> ▪ Using an AIC adaptive approach to select among a constant effects model, a linear interaction model, and a
- <sup>8</sup>     model using restricted cubic spline transformations was better suited for larger sample sizes.

9    **What this adds to what was known?**

- 10    ▪ Using a constant relative treatment effect for making risk-based individualized benefit predictions is not  
11    always the best approach.
- 12    ▪ Models including linear interaction of treatment with the prognostic index have robust performance across a  
13    wide variety of settings.
- 14    ▪ Increasing the flexibility of the transformation of the prognostic index does not improve performance.

15    **What is the implication and what should change now?**

- 16    ▪ Assuming a constant relative treatment effect for all patients may lead to treatment mistargeting.
- 17    ▪ We recommend using a linear interaction of treatment with the prognostic index for making individualized  
18    risk-based benefit predictions.
- 19    ▪ With larger sample size, using an adaptive approach based on AIC to explore non-linear interactions of the  
20    prognostic index with treatment is a viable option.

21    **1. Introduction**

22    Predictive approaches for assessing heterogeneity of treatment effects (HTE) aim at the development of  
23    models predicting either individualized effects or which of two (or more) treatments is better for an individual  
24    [1]. In prior work, we divided such methods in three broader categories based on the reference class used for  
25    defining patient similarity when making individualized predictions or recommendations [2]. First, risk-modeling  
26    approaches use prediction of baseline risk as the reference; second, treatment effect modeling approaches also  
27    model treatment-covariate interactions, in addition to risk factors; third, optimal treatment regime approaches  
28    focus on developing treatment assignment rules and rely heavily on modeling treatment effect modifiers. A key  
29    difference between these approaches is their parsimony in dealing the treatment effect modifiers, with no interaction  
30    considered (risk modeling), a limited number of interactions (effect modeling), or a larger set of interactions  
31    (optimal treatment regime approaches).

32    Risk-modeling approaches to predictive HTE analyses provide a viable option in the absence of well-established  
33    treatment effect modifiers [3,4]. In simulations, modeling of effect modifiers, i.e. treatment-covariate interactions,  
34    often led to miscalibrated predictions of absolute benefit, while risk-based methods proved quite robust in terms  
35    of benefit calibration, although provided weaker discrimination of benefit in the presence of true effect modifiers  
36    [5]. Most often, risk-modeling approaches are carried out in two steps: first a risk prediction model is developed  
37    externally or internally on the entire RCT population, “blinded” to treatment; then the RCT population is stratified  
38    using this prediction model to evaluate risk-based treatment effect variation [6]. This two-step approach identified

39 substantial absolute treatment effect differences between low-risk and high-risk patients in a re-analysis of 32 large  
 40 trials [7]. However, even though estimates at the risk subgroup level may be accurate, these estimates may need  
 41 further refinement for individual patients, especially for patients with predicted risk at the boundaries of the risk  
 42 intervals. Hence, the risk-stratified approach is useful for exploring and presenting HTE, but is not sufficient for  
 43 supporting treatment decisions for individual patients.

44 To individualize treatment effects, the recent PATH statement suggested various risk-based models including a  
 45 prognostic index of baseline risk (PI) and treatment assignment [3,4]. In the current simulation study, we aim to  
 46 summarize and compare different risk-based models for predicting individualized treatment effects. We simulate  
 47 different relations between baseline risk and treatment effects and also consider potential harms of treatment. We  
 48 illustrate the different models by a case study of predicting individualized effects of treatment for acute myocardial  
 49 infarction (MI) in a large randomized controlled trial (RCT).

## 50 2. Methods

### 51 2.1. Background

We observe RCT data  $(Z, X, Y)$ , where for each patient  $Z_i = 0, 1$  is the treatment status,  $Y_i = 0, 1$  is the observed outcome and  $X_i$  is a set of covariates measured. Let  $\{Y_i(z), z = 0, 1\}$  denote the unobservable potential outcomes. We observe  $Y_i = Z_i Y_i(1) + (1 - Z_i) Y_i(0)$ . We are interested in predicting the conditional average treatment effect (CATE),

$$\tau(x) = E\{Y(0) - Y(1)|X = x\}$$

Assuming that  $(Z, X, Y)$  is a random sample from the target population and that  $(Y(0), Y(1)) \perp\!\!\!\perp Z|X$ , as we are in the RCT setting, we can predict CATE from

$$\begin{aligned}\tau(x) &= E\{Y(0) | X = x\} - E\{Y(1) | X = x\} \\ &= E\{Y | X = x, Z = 0\} - E\{Y | X = x, Z = 1\}\end{aligned}$$

Based on an estimate of baseline risk

$$E\{Y | X = x, Z = 0\} = g(lp(x))$$

with  $u = lp(x) = \hat{\beta}^t x$  the linear predictor and  $g$  the link function, we predict CATE from

$$\hat{\tau}(x) = g(f(u, 0)) - g(f(u, 1))$$

52 where  $f(u, z)$  describes interactions of the baseline risk linear predictor with treatment. In the rest of the paper we  
53 use the linear predictor  $lp(x)$  as the prognostic index (PI) for the outcome  $Y$ .

54 *2.2. Simulation scenarios*

55 We simulated a typical RCT that is undertaken to compare a binary outcome (e.g. death) between a group of  
56 patients in the treatment arm and a group of untreated patients in the control arm. For each patient we generated  
57 8 baseline covariates  $x_1, \dots, x_4 \sim N(0, 1)$  and  $x_5, \dots, x_8 \sim B(1, 0.2)$ . Treatment was allocated using a 50:50  
58 split. Outcomes for patients in the control arm were generated from a logistic regression model including all  
59 baseline covariates. In the base scenarios coefficient values were such, that the AUC of the logistic regression  
60 model was 0.75 and the event rate in the control arm was 20%. Binary outcomes in the control arm were generated  
61 from Bernoulli variables with true probabilities  $P(y = 1|X, t_x = 0) = \text{expit}(PI) = \frac{e^{PI}}{1+e^{PI}}$ .

Outcomes in the treatment arm were generated using 3 base scenarios: absent treatment effect ( $OR = 1$ ),  
moderate treatment effect ( $OR = 0.8$ ) and strong treatment effect ( $OR = 0.5$ ). We started with simulating  
outcomes based on true constant relative treatment effects for the 3 base scenarios. We then simulated linear,  
quadratic and non-monotonic deviations from constant treatment effects using:

$$lp_1 = \gamma_2(PI - c)^2 + \gamma_1(PI - c) + \gamma_0,$$

62 where  $lp_1$  is the true linear predictor in the treatment arm, so that  $P(y = 1|X, Z = 1) = \text{expit}(lp_1)$ . Finally, we  
63 simulated scenarios where a constant absolute harm is applied across all treated patients. In this case we have  
64  $P(y = 1|X, Z = 1) = \text{expit}(lp_1) + \text{harm}$ .

65 The sample size for the base scenarios was set to 4,250, since this sample size provides 80% power for the  
66 detection of a marginal OR of 0.8 with the standard alpha of 0.5%. We evaluated the effect of smaller or larger  
67 sample sizes of 1,063 (4,250 divided by 4) and 17,000 (4250 multiplied by 4), respectively. We also evaluated the  
68 effect of worse or better discriminative ability for risk, adjusting the baseline covariate coefficients, such that the  
69 AUC of the regression model in the control arm was 0.65 and 0.85 respectively.

70 Combining all these settings resulted in a simulation study of 648 scenarios (exact settings in the supplementary  
71 material). With these scenarios we were able to cover the observed treatment effect heterogeneity in 32 large trials  
72 as well as many other potential variations of risk-based treatment effect [7].

73 *2.3. Individualized risk-based benefit predictions*

74 All risk-based methods assume that a risk prediction model is available to assign risk predictions to individual  
75 patients. For the simulations we developed a prediction model internally on the entire population, using a logistic  
76 regression model with main effects for all baseline covariates and treatment assignment. Baseline risk predictions

77 for individual patients were derived by setting treatment assignment to 0. Another common approach is to derive  
78 the prediction model solely on the control patients, however this approach has been shown to lead to biased benefit  
79 predictions [5,8,9].

80 A *stratified HTE method* has been suggested as an alternative to traditional subgroup analyses [3,4]. Patients  
81 are stratified into equally-sized risk strata—in this case based on risk quartiles. Absolute treatment effects within  
82 risk strata are estimated by the difference in event rate between patients in the control arm and patients in the  
83 treated arm. We considered this approach as a reference, expecting it to perform worse than the other candidates,  
84 as its objective is to provide an illustration of HTE rather than to optimize individualized benefit predictions.

85 Second, we considered a model which assumes *constant relative treatment effect* (constant odds ratio). Hence,  
86 absolute benefit is predicted from  $\hat{\tau}(\mathbf{x}) = \text{expit}(PI) - \text{expit}(PI + \log(OR))$ .

87 Third, we considered a logistic regression model including treatment, the prognostic index, and their linear  
88 interaction. Absolute benefit is then estimated from  $\hat{\tau}(\mathbf{x}) = \text{expit}(\beta_0 + \beta_{PI}PI) - \text{expit}(\beta_0 + \beta_{tx} + (\beta_{PI} + \beta_*)PI)$ .  
89 We will refer to this method as the *linear interaction* approach.

90 Fourth, we used *restricted cubic splines* (RCS) to relax the linearity assumption on the effect of the linear  
91 predictor [10]. We considered splines with 3 (RCS-3), 4 (RCS-4) and 5 (RCS-5) knots to compare models with  
92 different levels of flexibility.

93 Finally, we considered an *adaptive approach* using Akaike's Information Criterion (AIC) for model selection.  
94 More specifically, for the adaptive approach we ranked the constant relative treatment effect model, the linear  
95 interaction model, and the RCS models with 3, 4, and 5 knots based on their AIC and selected the one with  
96 the lowest value. The extra degrees of freedom were 1 (linear interaction), 2, 3 and 4 (RCS models) for these  
97 increasingly complex interactions with the treatment effect.

#### 98 2.4. Evaluation metrics

99 We evaluated the predictive accuracy of the considered methods by the root mean squared error (RMSE):

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\tau(\mathbf{x}_i) - \hat{\tau}(\mathbf{x}_i))^2}$$

100 We compared the discriminative ability of the methods under study using c-for-benefit [11]. The c-for-benefit  
101 represents the probability that from two randomly chosen matched patient pairs with unequal observed benefit,  
102 the pair with greater observed benefit also has a higher predicted benefit. To be able to calculate observed benefit,  
103 patients in each treatment arm are ranked based on their predicted benefit and then matched 1:1 across treatment  
104 arms. *Observed* treatment benefit is defined as the difference of observed outcomes between the untreated and  
105 the treated patient of each matched patient pair. *Predicted* benefit is defined as the average of predicted benefit

106 within each matched patient pair.

107 We evaluated calibration in a similar manner, using the integrated calibration index (ICI) for benefit [12]. The  
108 observed benefits are regressed on the predicted benefits using a locally weighted scatterplot smoother (loess).  
109 The ICI-for-benefit is the average absolute difference between predicted and smooth observed benefit. Values  
110 closer to 0 represent better calibration.

111 For each scenario setting we performed 500 replications, within which all the considered models were fitted.  
112 For comparing between models we simulated a super-population of size 500,000 for each scenario. We calculated  
113 RMSE and discrimination and calibration for benefit of the models derived in each replication of our simulation  
114 settings within this super-population.

115 *2.5. Empirical illustration*

116 We demonstrated the different methods for individualizing treatment benefits using data from 30,510 patients  
117 with acute myocardial infarction (MI) included in the GUSTO-I trial. 10,348 patients were randomized to tissue  
118 plasminogen activator (tPA) treatment and 20,162 were randomized to streptokinase. The outcome of interest  
119 was 30-day mortality (total of 2,128 events), recorded for all patients.

120 In line with previous analyses [13,14], we fitted a logistic regression model with 6 baseline covariates, i.e. age,  
121 Killip class, systolic blood pressure, heart rate, an indicator of previous MI, and the location of MI, to predict  
122 30-day mortality risk. A constant effect of treatment was included in the model. When deriving risk predictions  
123 for individuals we set the treatment indicator to 0. More information on model development can be found in the  
124 supplement (Supplement, Section 8).

125 **3. Results**

126 *3.1. Simulations*

127 The linear interaction model outperformed all RCS methods in terms of RMSE in scenarios with true constant  
128 relative treatment effect ( $OR = 0.8$ ,  $N = 4,250$  and  $AUC = 0.75$ ), strong linear and even strong quadratic deviations  
129 from a constant relative treatment effect (Figure 1; panels A-C). However, with non-monotonic deviations from a  
130 constant relative treatment effect, the RMSE of the linear interaction model increased substantially, especially in the  
131 presence of treatment-related harms (Figure 1; panel D). In these scenarios, RCS-3 outperformed all other methods  
132 in terms of RMSE. As might be expected the constant treatment effect approach had overall best performance  
133 under true constant treatment effect settings. It was sensitive to all considered deviations, resulting in increased  
134 RMSE. Finally, the adaptive approach had comparable performance to the best-performing method in each scenario.  
135 However, in comparison with the best-performing approach, its RMSE was more variable in the scenarios with  
136 linear and non-monotonic deviations, especially when also including moderate or strong treatment-related harms.

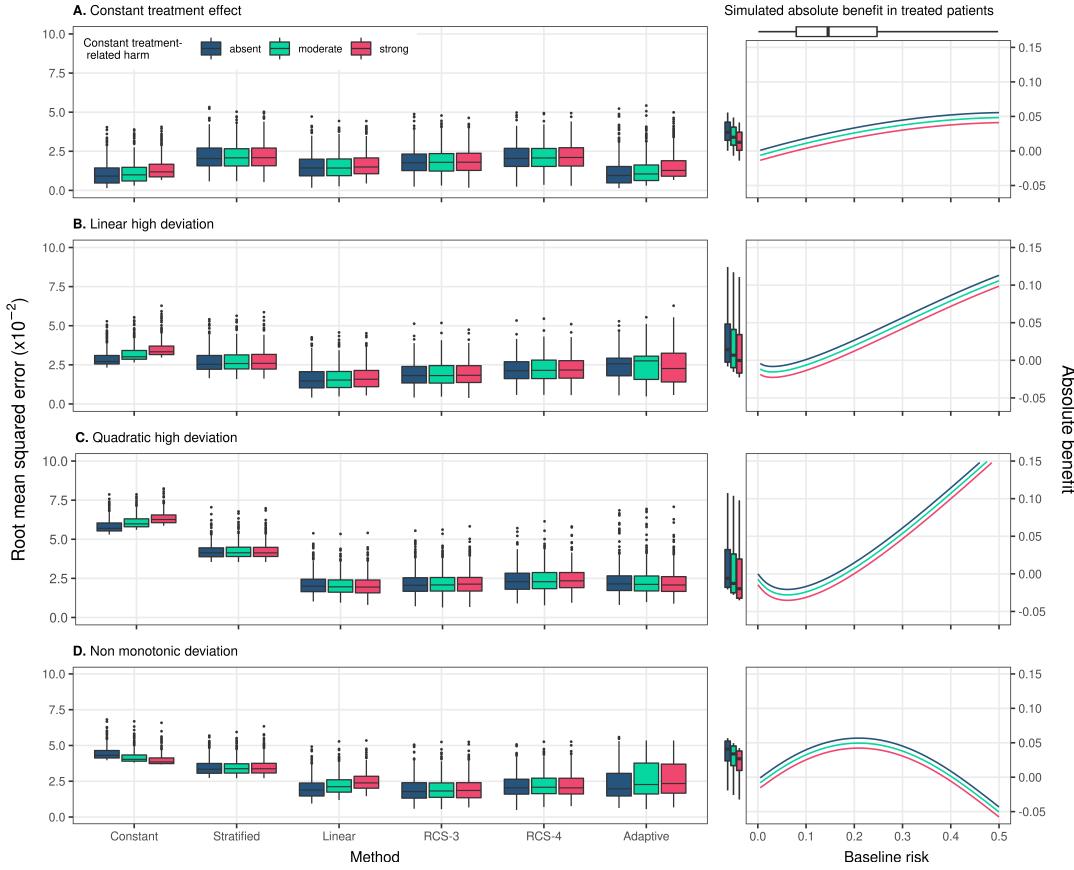


Figure 1: RMSE of the considered methods across 500 replications calculated from a simulated super-population of size 500,000. The scenario with true constant relative treatment effect (panel A) had a true prediction AUC of 0.75 and sample size of 4250. The RMSE is also presented for strong linear (panel B), strong quadratic (panel C), and non-monotonic (panel D) from constant relative treatment effects. Panels on the right side present the true relations between baseline risk (x-axis) and absolute treatment benefit (y-axis). The 2.5, 25, 50, 75, and 97.5 percentiles of the risk distribution are expressed by the boxplot on the top. The 2.5, 25, 50, 75, and 97.5 percentiles of the true benefit distributions are expressed by the boxplots on the side of the right-handside panel.

137 On closer inspection, we found that this behavior was caused by wrongly selecting the constant treatment effect  
 138 model in a substantial proportion of the replications (Supplement, Figure S3). This problematic behavior was less  
 139 with larger sample sizes (see below).

140 Increasing the sample size to 17,000 favored RCS-3 the most, It achieved lowest or close to lowest RMSE  
 141 across all scenarios (Figure 2). Especially in cases of strong quadratic and non-monotonic deviations RCS-3  
 142 had lower RMSE (median 0.011 for strong quadratic deviations and 0.010 for non-monotonic deviations with no  
 143 treatment-related harms) compared to the linear interaction approach (median 0.013 and 0.014, respectively),  
 144 regardless of the strength of treatment-related harms. Due to the large sample size, the RMSE of the adaptive  
 145 approach was even more similar to the best-performing method, and the constant relative treatment effect model  
 146 was less often wrongly selected (Supplement, Figure S4).

147 When we increased the AUC of the true prediction model to 0.85 (OR = 0.8 and N = 4,250). RCS-3 had the  
 148 lowest RMSE in the case of strong quadratic or non-monotonic deviations and very comparable performance to

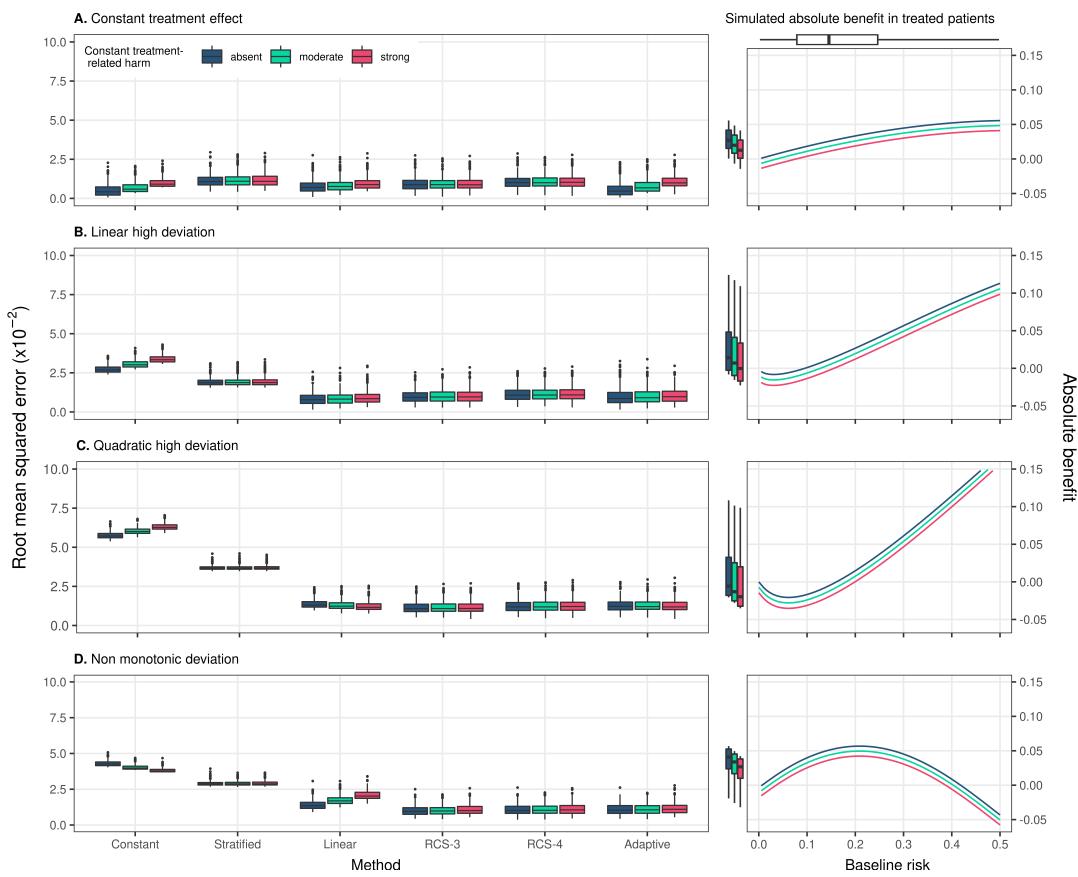


Figure 2: RMSE of the considered methods across 500 replications calculated in simulated samples of size 17,000 rather than 4,250 in Figure 1. RMSE was calculated on a super-population of size 500,000

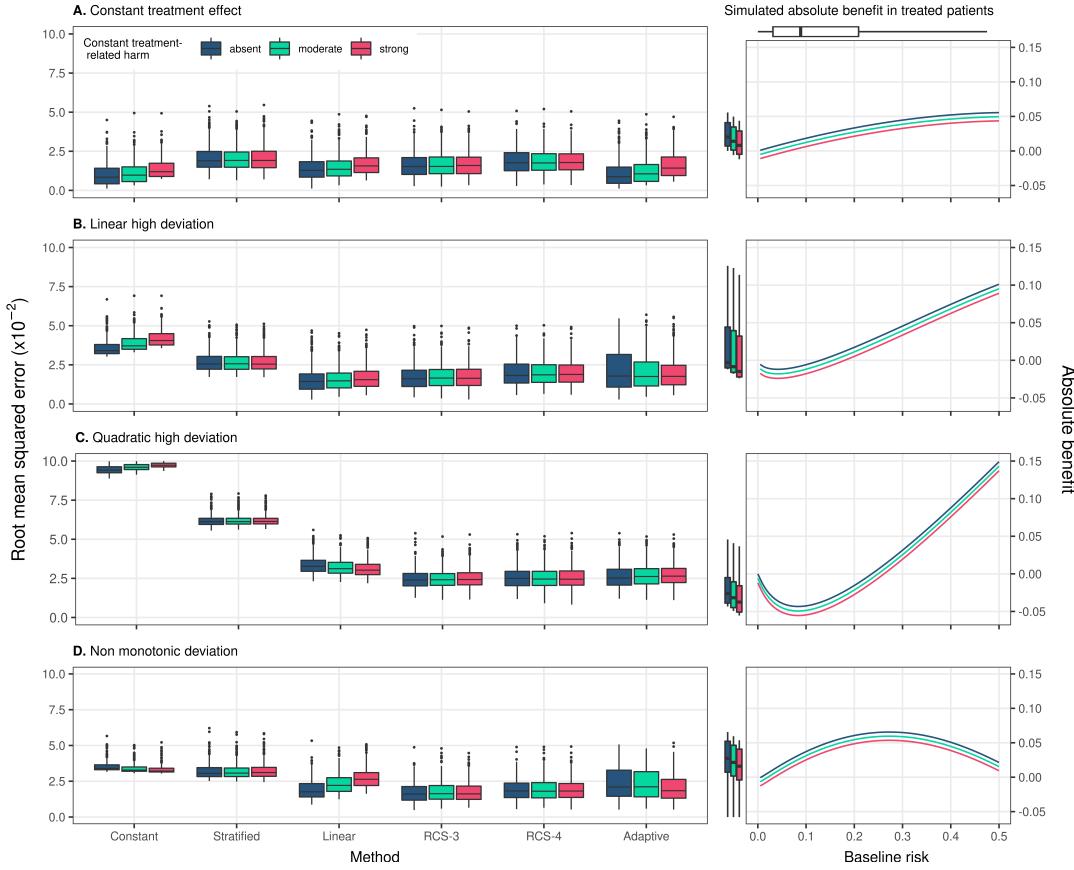


Figure 3: RMSE of the considered methods across 500 replications calculated in simulated samples 4,250. True prediction AUC of 0.85. RMSE was calculated on a super-population of size 500,000

149 the – optimal – linear interaction model in the case of strong linear deviations (median RMSE 0.016 for RCS-3  
 150 compared to 0.014 for the linear interaction model). As observed in the base case scenario the adaptive approach  
 151 wrongly selected the constant treatment effect model (23% and 25% of the replications in the strong linear and  
 152 non-monotonic deviation scenarios without treatment-related harms, respectively), leading to more variability of  
 153 the RMSE (Supplement, Figure S5).

154 When assuming a true constant relative treatment effect, discrimination for benefit was only slightly lower  
 155 for the linear interaction model, but substantially lower for the non-linear RCS approaches (Figure 4; panel A).  
 156 With strong linear or quadratic deviations from a constant relative treatment effect, all methods discriminated  
 157 quite similarly (Figure 4; panels B-C). In the scenario with non-monotonic deviations, the constant effect model  
 158 had much lower discriminative ability compared to all other methods (median AUC of 0.4971 for the constant  
 159 effects model, 0.5285 for the linear interaction model and 0.5304 for the best-performing RCS-3; Figure 4; panel  
 160 D). The adaptive approach was unstable in terms of discrimination for benefit, especially in the presence of  
 161 treatment-related harms. With increasing number of RCS knots, we observed decreasing median values and  
 162 increasing variability of the c-for-benefit in all scenarios. When we increased the sample size to 17,000 we observed

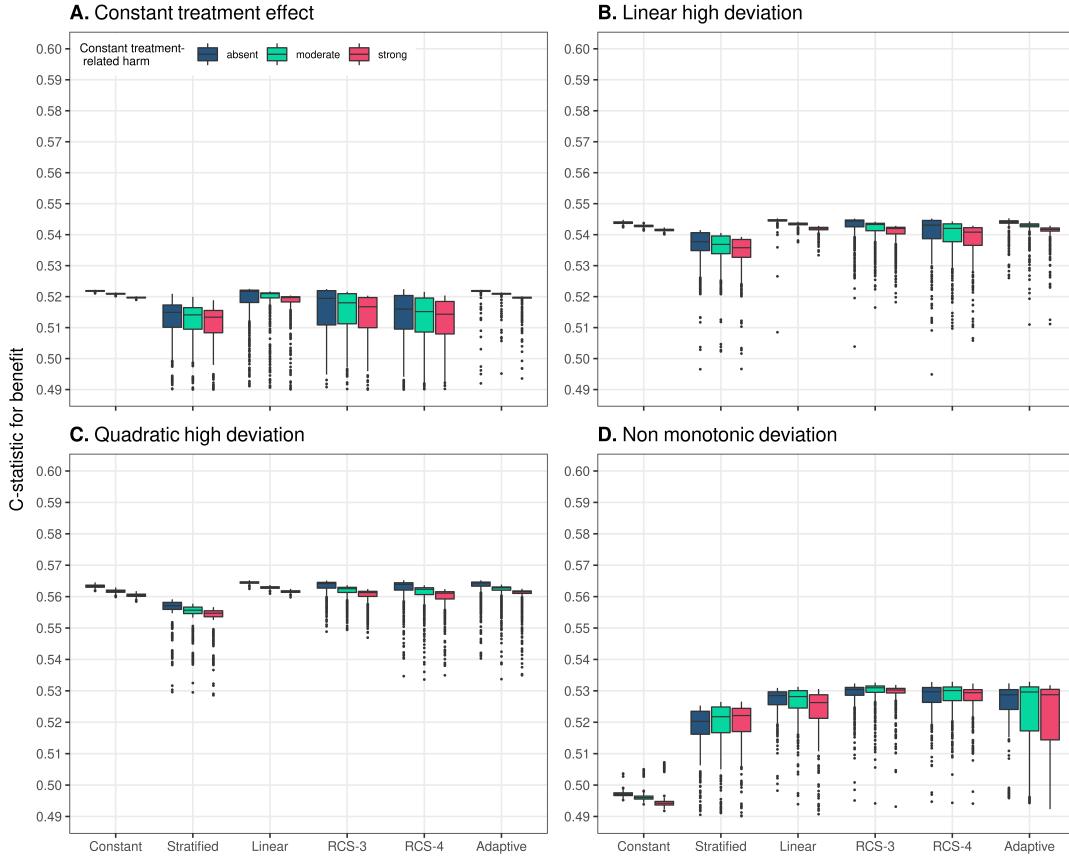


Figure 4: Discrimination for benefit of the considered methods across 500 replications calculated in a simulated samples of size 4,250. True prediction AUC of 0.75.

163 similar trends, however the performance of all methods was more stable (Supplement, Figure S6). Finally, when  
 164 we increased the true prediction AUC to 0.85 the adaptive approach in the case of non-monotonic deviations was,  
 165 again, more conservative, especially with null or moderate treatment-related harms (Supplement, Figure S5).

166 In terms of calibration for benefit, the constant effects model outperformed all other models in the scenario with  
 167 true constant treatment effects, but was miscalibrated for all deviation scenarios (Figure 5). The linear interaction  
 168 model showed best or close to best calibration across all scenarios and only showed worse calibration compared  
 169 to RCS-3 in case of non-monotonic deviations and treatment-related harms (Figure 5; panel D). The adaptive  
 170 approach was worse calibrated in scenarios with strong linear and non-monotonic deviations compared to the linear  
 171 interaction model and RCS-3. When we increased sample size to 17,000 similar conclusions on calibration for  
 172 benefit could be drawn. As expected, all methods displayed more stable calibration performance due to the larger  
 173 number of patients (Supplement, Figure S6). When we increased the true prediction AUC to 0.85, the linear  
 174 interaction model was worse calibrated, on average, than RCS-3 in the case of strong quadratic deviations from  
 175 constant relative treatment effects (Supplement, Figure S7).

176 The results from all individual scenarios can be explored online at [https://arekkas.shinyapps.io/simulation\\_](https://arekkas.shinyapps.io/simulation_)

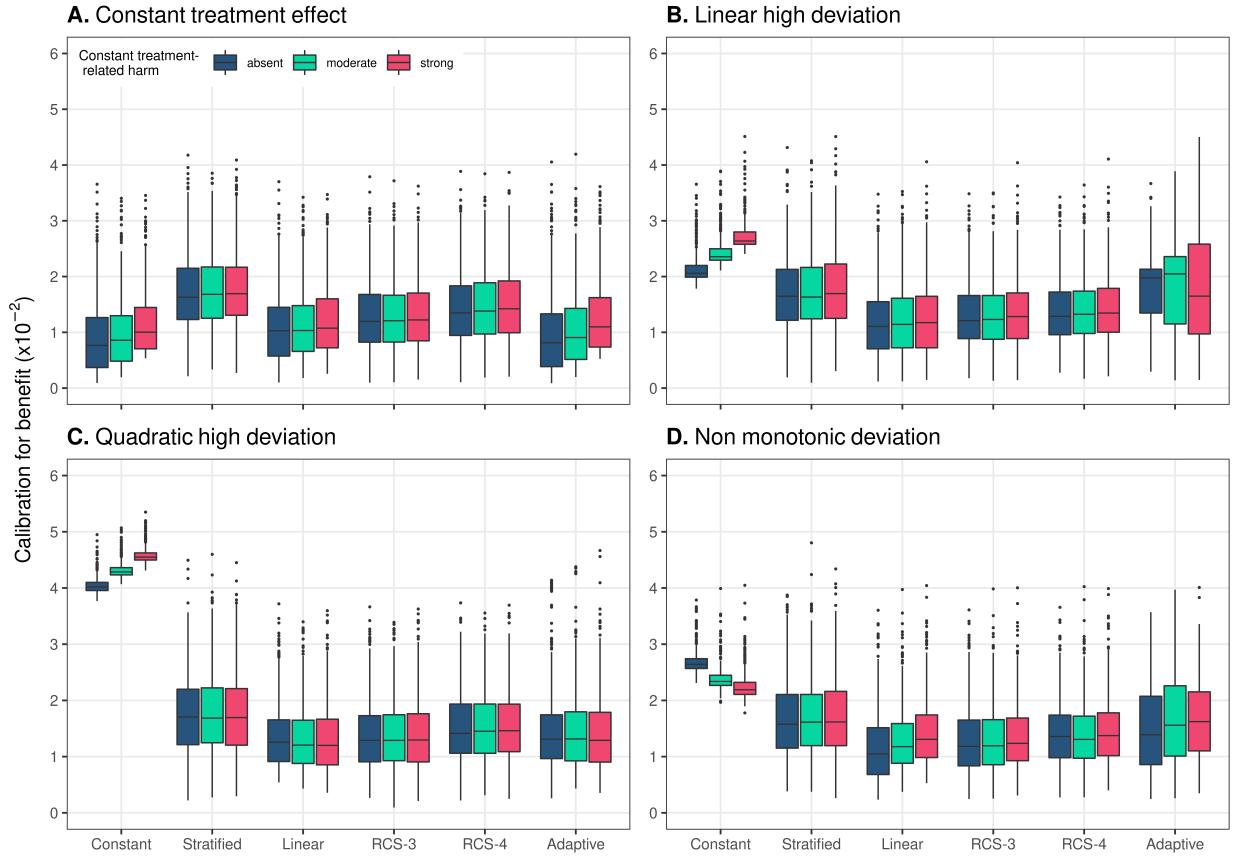


Figure 5: Calibration for benefit of the considered methods across 500 replications calculated in a simulated sample of size 500,000. True prediction AUC of 0.75 and sample size of 4,250.

177 viewer/. Additionally, all the code for the simulations can be found at [https://github.com/rekkasa/arekkas\\_](https://github.com/rekkasa/arekkas_)  
 178 HteSimulation\_XXXX\_2021

### 179 3.2. Empirical illustration

180 We used the derived prognostic index to fit a constant treatment effect model, a linear interaction model  
 181 and a RCS-3 model individualizing absolute benefit predictions. RCS-4 and RCS-5 models were excluded. In our  
 182 simulations these methods were always outperformed by the simpler approaches and were often overfitted. Finally,  
 183 an adaptive approach with only the 3 candidate models was also applied.

184 All considered methods provided similar fits, predicting increasing benefits for patients with higher baseline risk  
 185 predictions. All models followed the evolution of the stratified estimates very closely. The adaptive approach based  
 186 on AIC selected the constant treatment effect model. The constant treatment effect model had somewhat lower  
 187 AIC compared to the linear interaction model slightly worse cross-validated discrimination (c-for-benefit 0.525 vs  
 188 0.526) and better cross-validated calibration (ICI-for benefit 0.0104 vs 0.0115). In conclusion, a simpler constant  
 189 treatment effect model is adequate for predicting absolute 30-day mortality benefits of treatment with tPA in

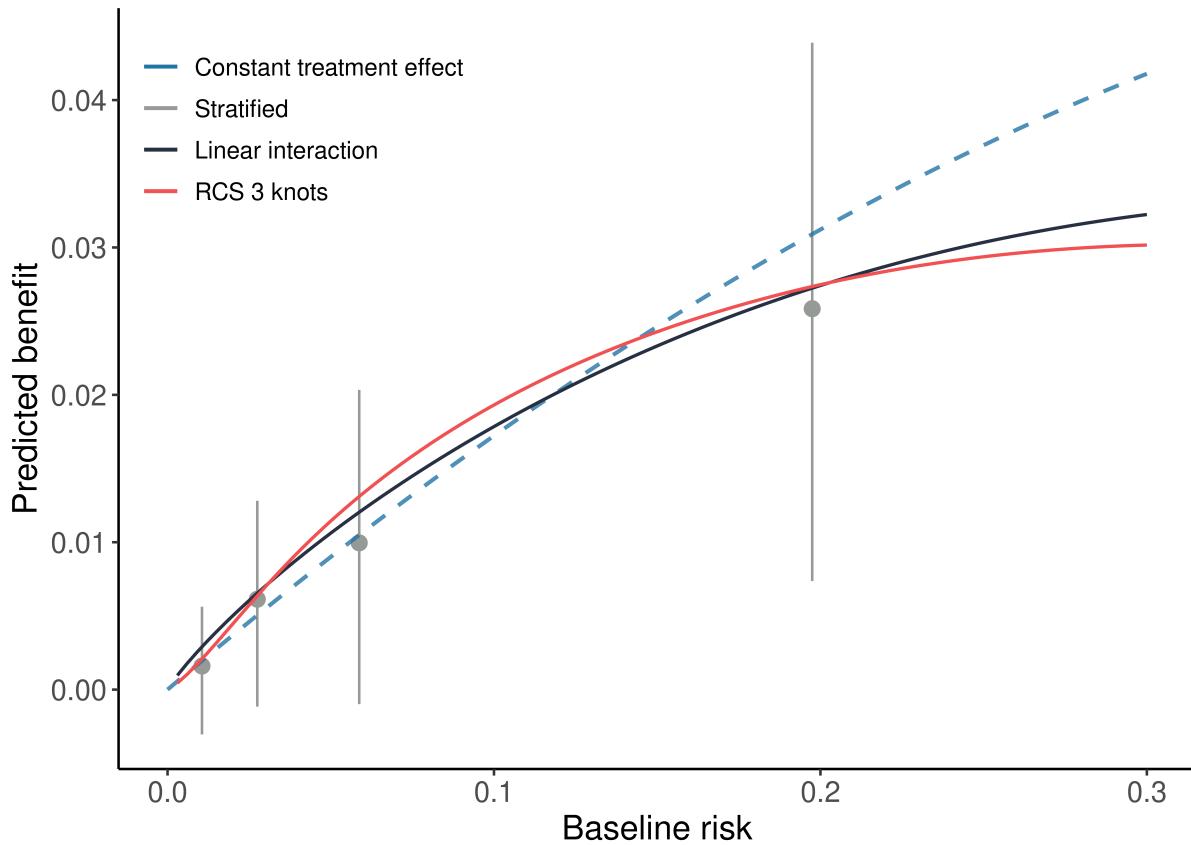


Figure 6: Individualized absolute benefit predictions based on baseline risk when using a constant treatment effect approach, a linear interaction approach and RCS smoothing using 3 knots. Risk stratified estimates of absolute benefit are presented within quartiles of baseline risk as reference.

<sup>190</sup> patients with acute MI.

#### <sup>191</sup> 4. Discussion

<sup>192</sup> The linear interaction model and the RCS-3 model both displayed very good performance under many of the  
<sup>193</sup> considered simulation scenarios. The linear interaction model was optimal in cases with smaller sample sizes and  
<sup>194</sup> moderately performing baseline risk prediction models, that is, it had lower RMSE, was better calibrated for benefit  
<sup>195</sup> and had better discrimination for benefit, even in scenarios with strong quadratic deviations. In scenarios with true  
<sup>196</sup> non-monotonic deviations, the linear interaction model was outperformed by RCS-3, especially in the presence of  
<sup>197</sup> true treatment-related harms. Increasing the sample size or the prediction model's discriminative ability favored  
<sup>198</sup> RCS-3, especially in scenarios with non-monotonic deviations and in the presence of treatment-related harms.

<sup>199</sup> RCS-4 and RCS-5 proved to be too flexible in all considered scenarios, as indicated by higher RMSE, increased  
<sup>200</sup> variability of discrimination for benefit and worse calibration of benefit predictions. Even with larger sample sizes  
<sup>201</sup> and strong quadratic or non-monotonic deviations from the base case scenario of constant relative treatment

202 effects, these more flexible restricted cubic splines did not outperform the simpler RCS-3. These approaches may  
203 only be helpful if we expect more extreme patterns of heterogeneous treatment effects compared to the quadratic  
204 deviations considered here. Considering interactions in RCS-3 models as the most complex approach often may be  
205 reasonable.

206 The constant treatment effect model, despite having adequate performance in the presence of weak treatment  
207 effect heterogeneity on the relative scale, quickly broke down with stronger deviations from constant relative  
208 treatment effects. In these cases, the stratified approach generally had lower error rates compared to the constant  
209 treatment effect model. Such stepwise treatment benefit estimates are useful for visually demonstrating treatment  
210 effect heterogeneity but may be considered insufficient for making individualized benefit predictions.

211 Increasing the discriminative ability of the risk model—by increasing the predictor coefficients of the true risk  
212 model—reduced RMSE for all methods. This increase in discriminative ability translates in higher variability of  
213 predicted risks, which, in turn, allows the considered methods to better capture absolute treatment benefits. As  
214 a consequence, the increase in discriminative ability of the risk model also led to higher discrimination between  
215 those with low or high benefit (as reflected in values of c-for-benefit). Even though risk model performance is very  
216 important for the ability of risk-based methods to predict treatment benefit, prediction model development was  
217 outside the scope of this work and has already been studied extensively [5,8,9].

218 The adaptive approach had adequate performance, following closely on average the performance of the “true”  
219 model in most scenarios. With smaller sample sizes it tended to miss the treatment-risk interactions and selected  
220 simpler models (Supplement Section 4). This conservative behavior resulted in increased RMSE variability in these  
221 scenarios, especially in the case of true strong linear or non-monotonic deviations from the base case scenario.  
222 Therefore, in the case of smaller sample sizes the simpler linear interaction model may be a safer choice for  
223 predicting absolute benefits in the presence of any suspected treatment-related harms.

224 A limitation of our study is that we assumed treatment benefit to be a function of baseline risk in the majority  
225 of the simulation scenarios. We also considered constant moderate and strong treatment-related harms, applied on  
226 the absolute scale to expand the range of scenarios in line with previous work [15]. In a limited set of scenarios  
227 where we assumed the existence of true treatment-covariate interactions, our conclusions remained unchanged.  
228 Even though the average error rates increased for all the considered methods, due to the miss-specification of the  
229 outcome model, the linear interaction model had the lowest error rates. RCS-3 had very comparable performance.  
230 The constant treatment effect model often gave biased results, especially in the presence of moderate or strong  
231 treatment-related harms. All the results of these simulations can be found in Supplement, Section 7. Future  
232 simulation studies could explore the effect of more extensive deviations from risk-based treatment effects.

233 In our simulations we only focused on risk-based methods, using baseline risk as a reference in a two-stage  
234 approach to individualizing benefit predictions. However, there is a plethora of different methods, ranging from

235 treatment effect modeling to tree-based approaches available in more recent literature [16–19]. Many of these  
236 methods rely on incorporating treatment-covariate interactions in the prediction of benefit. An important caveat  
237 of such approaches is that they may be prone to overfitting, thus exaggerating the magnitude of the predicted  
238 benefits. In a wide range of simulation settings, a simpler risk modeling approach was consistently better calibrated  
239 for benefit compared to more complex treatment effect modelling approaches [5]. However, whether this remains  
240 the case in a range of empirical settings still needs to be explored. Similarly, when SYNTAX score II, a model  
241 developed for identifying patients with complex coronary artery disease that benefit more from percutaneous  
242 coronary intervention or from coronary artery bypass grafting was redeveloped using fewer treatment-covariate  
243 interactions had better external performance compared to its predecessor[20,21].

244 In conclusion, the linear interaction approach is a viable option with smaller sample sizes and/or moderately  
245 performing risk prediction models if we consider a non-constant relative treatment effect plausible. RCS-3 is  
246 a better option when non-monotonic deviations from a constant relative treatment effect and/or substantial  
247 treatment-related harms are anticipated. Increasing the complexity of the RCS models by increasing the number  
248 of knots does not translate to improved benefit prediction. Using AIC for model selection among the constant  
249 treatment effect, the linear interaction and RCS-3 model is a viable option, especially with larger sample size.

250 **5. References**

- 251 [1] Varadhan R, Segal JB, Boyd CM, Wu AW, Weiss CO. A framework for the analysis of heterogeneity of  
252 treatment effect in patient-centered outcomes research. *Journal of Clinical Epidemiology* 2013;66:818–25.  
253 <https://doi.org/10.1016/j.jclinepi.2013.02.009>.
- 254 [2] Rekkas A, Paulus JK, Raman G, Wong JB, Steyerberg EW, Rijnbeek PR, et al. Predictive approaches  
255 to heterogeneous treatment effects: A scoping review. *BMC Medical Research Methodology* 2020;20.  
256 <https://doi.org/10.1186/s12874-020-01145-1>.
- 257 [3] Kent DM, Paulus JK, Klaveren D van, D'Agostino R, Goodman S, Hayward R, et al. The predictive  
258 approaches to treatment effect heterogeneity (PATH) statement. *Annals of Internal Medicine* 2019;172:35.  
259 <https://doi.org/10.7326/m18-3667>.
- 260 [4] Kent DM, Klaveren D van, Paulus JK, D'Agostino R, Goodman S, Hayward R, et al. The predictive approaches  
261 to treatment effect heterogeneity (PATH) statement: Explanation and elaboration. *Annals of Internal Medicine*  
262 2019;172:W1. <https://doi.org/10.7326/m18-3668>.
- 263 [5] Klaveren D van, Balan TA, Steyerberg EW, Kent DM. Models with interactions overestimated heterogeneity of  
264 treatment effects and were prone to treatment mistargeting. *Journal of Clinical Epidemiology* 2019;114:72–83.  
265 <https://doi.org/10.1016/j.jclinepi.2019.05.029>.
- 266 [6] Kent DM, Rothwell PM, Ioannidis JP, Altman DG, Hayward RA. Assessing and reporting heterogeneity in  
267 treatment effects in clinical trials: A proposal. *Trials* 2010;11. <https://doi.org/10.1186/1745-6215-11-85>.
- 268 [7] Kent DM, Nelson J, Dahabreh IJ, Rothwell PM, Altman DG, Hayward RA. Risk and treatment effect  
269 heterogeneity: Re-analysis of individual participant data from 32 large clinical trials. *International Journal of  
270 Epidemiology* 2016;dyw118. <https://doi.org/10.1093/ije/dyw118>.
- 271 [8] Burke JF, Hayward RA, Nelson JP, Kent DM. Using internally developed risk models to assess heterogeneity  
272 in treatment effects in clinical trials. *Circulation: Cardiovascular Quality and Outcomes* 2014;7:163–9.  
273 <https://doi.org/10.1161/circoutcomes.113.000497>.
- 274 [9] Abadie A, Chingos MM, West MR. Endogenous stratification in randomized experiments. *The Review of  
275 Economics and Statistics* 2018;100:567–80. [https://doi.org/10.1162/rest\\_a\\_00732](https://doi.org/10.1162/rest_a_00732).
- 276 [10] Harrell FE, Lee KL, Pollock BG. Regression models in clinical studies: Determining relationships between  
277 predictors and response. *JNCI Journal of the National Cancer Institute* 1988;80:1198–202. <https://doi.org/10.1093/jnci/80.15.1198>.
- 279 [11] Klaveren D van, Steyerberg EW, Serruys PW, Kent DM. The proposed “concordance-statistic for benefit”  
280 provided a useful metric when modeling heterogeneous treatment effects. *Journal of Clinical Epidemiology*  
281 2018;94:59–68. <https://doi.org/10.1016/j.jclinepi.2017.10.021>.

- 282 [12] Austin PC, Steyerberg EW. The integrated calibration index (ICI) and related metrics for quantifying the  
283 calibration of logistic regression models. *Statistics in Medicine* 2019;38:4051–65. <https://doi.org/10.1002/sim.8281>.
- 285 [13] Califf RM, Woodlief LH, Harrell FE, Lee KL, White HD, Guerci A, et al. Selection of thrombolytic  
286 therapy for individual patients: Development of a clinical model. *American Heart Journal* 1997;133:630–9.  
287 [https://doi.org/10.1016/s0002-8703\(97\)70164-9](https://doi.org/10.1016/s0002-8703(97)70164-9).
- 288 [14] Steyerberg EW, Bossuyt PMM, Lee KL. Clinical trials in acute myocardial infarction: Should we adjust  
289 for baseline characteristics? *American Heart Journal* 2000;139:745–51. [https://doi.org/10.1016/s0002-8703\(00\)90001-2](https://doi.org/10.1016/s0002-8703(00)90001-2).
- 291 [15] Glasziou PP, Irwig LM. An evidence based approach to individualising treatment. *BMJ* 1995;311:1356–9.  
292 <https://doi.org/10.1136/bmj.311.7016.1356>.
- 293 [16] Athey S, Tibshirani J, Wager S. Generalized random forests. *The Annals of Statistics* 2019;47. <https://doi.org/10.1214/18-aos1709>.
- 295 [17] Lu M, Sadiq S, Feaster DJ, Ishwaran H. Estimating individual treatment effect in observational data  
296 using random forest methods. *Journal of Computational and Graphical Statistics* 2018;27:209–19. <https://doi.org/10.1080/10618600.2017.1356325>.
- 298 [18] Wager S, Athey S. Estimation and inference of heterogeneous treatment effects using random forests. *Journal*  
299 *of the American Statistical Association* 2018;113:1228–42. <https://doi.org/10.1080/01621459.2017.1319839>.
- 301 [19] Powers S, Qian J, Jung K, Schuler A, Shah NH, Hastie T, et al. Some methods for heterogeneous treatment  
302 effect estimation in high dimensions. *Statistics in Medicine* 2018;37:1767–87.
- 304 [20] Farooq V, Van Klaveren D, Steyerberg EW, Meliga E, Vergouwe Y, Chieffo A, et al. Anatomical and  
305 clinical characteristics to guide decision making between coronary artery bypass surgery and percutaneous  
306 coronary intervention for individual patients: Development and validation of syntax score ii. *The Lancet*  
2013;381:639–50.
- 308 [21] Takahashi K, Serruys PW, Fuster V, Farkouh ME, Spertus JA, Cohen DJ, et al. Redevelopment and validation  
309 of the syntax score ii to individualise decision making between percutaneous and surgical revascularisation in  
patients with complex coronary artery disease: Secondary analysis of the multicentre randomised controlled  
syntaxes trial with external cohort validation. *The Lancet* 2020;396:1399–412.