

## **Performance metrics for models designed to predict treatment effect**

C.C.H.M. Maas<sup>\*1</sup>, D.M. Kent<sup>2</sup>, M.C. Hughes<sup>2</sup>, R. Dekker<sup>3</sup>, H.F. Lingsma<sup>1</sup>, D. van Klaveren<sup>1, 2</sup>

<sup>1</sup>Department of Public Health, Erasmus University Medical Center, Rotterdam, Netherlands

<sup>2</sup>Predictive Analytics and Comparative Effectiveness Center, Institute for Clinical Research and Health Policy Studies, Tufts Medical Center, Boston, USA

<sup>3</sup>Econometric Institute, Erasmus University Rotterdam, Rotterdam, Netherlands

### **Materials & Correspondence**

C.C.H.M. Maas, Erasmus University Medical Center, Doctor Molewaterplein 40, 3015 GD Rotterdam, Netherlands, [c.h.m.maas@erasmusmc.nl](mailto:c.h.m.maas@erasmusmc.nl)

## **ABSTRACT**

### **Objective**

Measuring the performance of models that predict individualized treatment effect is challenging because the outcomes of two alternative treatments are inherently unobservable in one patient. The C-for-benefit was proposed to measure discriminative ability. We aimed to propose metrics of calibration and overall performance for models predicting treatment effect.

### **Study Design and Setting**

Similar to the previously proposed C-for-benefit, we defined observed pairwise treatment effect as the difference between outcomes in pairs of matched patients with different treatment assignment. We redefined the E-statistics, the cross-entropy, and the Brier score into metrics for measuring a model's ability to predict treatment effect. In a simulation study, the metric values of deliberately "perturbed models" were compared to those of the data-generating model, i.e., "optimal model". To illustrate the performance metrics, models predicting treatment effects were applied to the data of the Diabetes Prevention Program.

### **Results**

As desired, performance metric values of "perturbed models" were consistently worse than those of the "optimal model" ( $E_{\text{avg-for-benefit}} \geq 0.070$  versus 0.001,  $E_{90\text{-for-benefit}} \geq 0.115$  versus 0.003, cross-entropy-for-benefit  $\geq 0.757$  versus 0.733, Brier-for-benefit  $\geq 0.215$  versus 0.212). Calibration, discriminative ability, and overall performance of three different models were similar in the case study.

## **Conclusion**

The proposed metrics are useful to assess the calibration and overall performance of models predicting treatment effect.

**Keywords:** heterogeneous treatment effect, prediction models, logistic regression, causal forest

## ARTICLE

### 1. Introduction

Clinicians and patients generally select the treatment that is expected to be beneficial on average for the patient population. However, the average treatment effect (ATE) for a population does not accurately reflect the effect of treatment for each patient individually[1-3]. Various models have been proposed for predicting individualized treatment effects[4]. These models aim to predict the difference between the outcomes of two alternative treatments for each patient.

Usually, only one of the outcomes can be observed for a given patient, the counterfactual outcome remains unobserved. This phenomenon—known as the fundamental problem of causal inference—complicates the assessment of a model's ability to predict treatment effect. As a result, the performance of models that predict treatment effect cannot be quantified with conventional metrics evaluating risk predictions[5]. To resolve this issue, observed pairwise treatment effect can be defined as the difference between outcomes in pairs of matched patients. Recently, the C-for-benefit has been proposed for quantifying to what extent the models can discriminate between patients who benefit and those who do not[6]. However, measures of calibration—the agreement between predicted and observed treatment effect in *groups* of patients—and measures of overall performance—the discrepancy between predicted and observed treatment effect across *individual* patients—are still lacking.

For models predicting outcome risk and not treatment effect, several metrics are available to assess calibration (i.e., E-statistic), and overall performance (i.e., cross-entropy and Brier score)[7-9]. However, these metrics may poorly reflect a model's ability to predict treatment effect. For example, in a simulation scenario with

a relatively small simulated data sample, the risk predictions of a model with all possible treatment interactions are reasonably well calibrated (Figure 1A), while the corresponding treatment effect predictions are poorly calibrated (Figure 1B)[10]. Apart from such graphical assessment of calibration in groups of patients with similar predicted treatment effects, no metrics are available that quantify the calibration or the overall performance of treatment effect predictions[11].

Therefore, we aimed to extend these performance metrics for calibration and overall performance for risk prediction models to models that are designed to predict treatment effect.

## 2. Methods

### 2.1 Definition of treatment effect

With the potential outcomes framework, we can define the (conditional average) treatment effect  $\tau(x)$  for a patient with baseline characteristics  $X$  as the expected difference between the outcome under control treatment  $Y_i(0)$  and the outcome under active treatment  $Y_i(1)$ , conditional on the patient characteristics  $X$ , i.e.

$$\tau(x) = E[Y_i(0) - Y_i(1)|X_i = x],$$

where the potential outcomes  $Y_i(W_i)$  indicate the outcome  $Y_i$  conditional on treatment  $W_i$ [12]. Here, the event associated with the outcome was assumed to be unfavorable. Thus, treatment benefit, i.e., a positive  $\tau(x)$ , is expected when the outcome probability under control treatment is higher than the outcome probability under active treatment. Alternatively, two active treatments can be administered instead of control and active treatment.

## **2.2 Metrics based on the matching principle**

Using the matching principle, we defined observed pairwise treatment effect as the difference in outcomes between two similar patients with different treatment assignments (Supplementary Information 1)[6]. Similarity was based on baseline patient characteristics to create pairs of similar patients with different treatment assignments. Specifically, we matched each untreated patient with the nearest treated patient based on the Mahalanobis distance between the patient characteristics without replacement[13]. With a binary outcome (say, 0 for alive and 1 for dead), four outcome combinations are possible for a pair of patients. First, treatment benefit was indicated if the treated patient lives and the untreated patient dies. Second, treatment harm was indicated if the treated patient dies and the untreated patient lives. Lastly, no effect of treatment was indicated if both the treated and untreated patients live, or if both die. Thus, the observed pairwise treatment effect takes the values 1 (benefit), 0 (no effect), and -1 (harm). Concurrently, predicted treatment effect is the difference between the predicted outcome probability of the untreated patient minus the predicted outcome probability of the treated patient (Supplementary Information 1). All of the following metrics use this matching principle and are added to Figure 1C for illustration.

### **2.2.1 Calibration**

Calibration refers to the correspondence between the predicted and observed treatment effects. The calibration-in-the-large or mean calibration was defined as the average observed pairwise treatment effect minus the average predicted treatment effect[14]. If the algorithm overestimates treatment effect, the average predicted treatment effect is higher than the observed pairwise treatment effect, resulting in a

negative calibration-in-the-large value. Conversely, the calibration-in-the-large will be positive if treatment effect is underestimated.

Calibration can also be assessed by a smoothed calibration curve obtained by a local regression of observed pairwise treatment effect on predicted treatment effect, with default values for the span and the degree of polynomials (Figure 1C). Similar to the E-statistic and the Integrated Calibration Index (ICI), we propose to measure calibration by the average absolute vertical distance between this smoothed calibration curve and the diagonal line of perfect calibration[7]. This quantity, which we named the  $E_{\text{avg}}$ -for-benefit, can be interpreted as the weighted difference between observed pairwise treatment effect and predicted treatment effect, with weights determined by the empirical density function of the predicted treatment effect. Similarly, we defined the  $E_{50}$ -for-benefit and the  $E_{90}$ -for-benefit as the median and 90<sup>th</sup> percentile of the absolute differences between the predicted treatment effect and the smoothed observed pairwise treatment effect (Supplementary Information 1)[7]. Thus, the E-statistics indicate perfect calibration when zero.

### **2.2.2 Discrimination**

Discrimination refers to a model's ability to separate patients who benefit from treatment and those who do not. To measure discrimination, we used the previously proposed C-for-benefit, i.e., the probability that from two randomly chosen matched patient pairs with unequal observed pairwise treatment effect, the pair with greater observed pairwise treatment effect also has a larger predicted treatment effect[6]. The C-for-benefit was calculated by the number of concordant pairs divided by the number of concordant and discordant pairs. Two patient pairs are concordant if the pair with the larger observed pairwise treatment effect also has a larger predicted

treatment effect. Two patient pairs are discordant if the pair with larger observed benefit has a smaller predicted treatment effect. Two patient pairs are uninformative if the pairs have the same observed pairwise treatment effect. The C-for-benefit is 0.5 if the model cannot distinguish between patients any better than random treatment assignment, and 1 if the model can perfectly distinguish between patients who benefit from treatment and who do not.

### **2.2.3 Overall performance measures**

We propose to measure overall performance using the multi-class versions of the Brier score and cross-entropy because observed pairwise treatment effect can belong to one of three classes (benefit, no effect, harm)[8, 9]. We defined cross-entropy-for-benefit as the logarithmic distance between predicted and observed pairwise treatment effect and Brier-for-benefit as the average squared distance between predicted and observed pairwise treatment effect (Supplementary Information 2). Thus, the overall performance metrics indicate better optimal performance when closer to zero. The cross-entropy-for-benefit and Brier-for-benefit measure overall model performance since these metrics are affected by calibration and discrimination simultaneously. The proposed metrics were implemented in a publicly available R-package “HTEPredictionMetrics”[15].

### **2.3 Data**

To illustrate the proposed metrics, we used data from the Diabetes Prevention Program (DPP). The participants of DPP were at risk to develop diabetes, which is defined as a body mass index of 24 or higher and impaired glucose metabolism[16]. The participants were randomized between 1996 and 2001 to receive 1) an intensive



program of lifestyle modification lessons, 2) 850 mg of metformin twice a day and standard lifestyle modification, or 3) placebo twice a day and standard lifestyle recommendations. To predict the effect of the intervention on the outcome, i.e., the risk of developing diabetes, we used the patient characteristics sex, age, ethnicity, body mass index, smoking status, fasting blood sugar, triglycerides, hemoglobin, self-reported history of hypertension, family history of diabetes, self-reported history of high blood glucose, and gestational diabetes mellitus (Supplementary Table 1). We imputed missing values of patient characteristics using Multivariate Imputations by Chained Equations (MICE)[17].

## 2.4 Simulation study

We simulated the outcomes of the DPP using the patient characteristics to study if the proposed performance metrics were better for the model used for outcome generation (“optimal model”) than for deliberately “perturbed models”. The “optimal model” was a logistic regression to model the probability of the outcome (developing diabetes)  $p_i$  based on the treatment (e.g., lifestyle intervention) assignment indicator  $W$ , a centered prognostic index  $PI$ , and their interaction:

$$\log \frac{p_i}{1 - p_i} = W_i \cdot \beta_W + (1 - W_i) \cdot PI_i \cdot \beta_{(1-W) \cdot PI} + W_i \cdot PI_i \cdot \beta_{W \cdot PI}.$$

The prognostic index  $PI$  ( $= X' \hat{\beta}_X$ ) was determined by regressing the outcome variable on the patient characteristics.

Next, we created a super population by duplicating the matched patient pairs 500 times to obtain high precision to ensure that observed differences between metrics are “true” differences. The outcomes of the super population  $Y_i$  were simulated from a Bernoulli distribution with the outcome probabilities  $p_i$  generated by the “optimal model”.

We then created three deliberate perturbations of the “optimal model”. The first “perturbed model” overestimates ATE by multiplying the coefficient of the treatment assignment indicator ( $\beta_W$ ) with 2 (Supplementary Figure 1). The second “perturbed model” overestimates risk heterogeneity by multiplying the coefficient of the prognostic index for both control treatment ( $\beta_{(1-W) \cdot PI}$ ) and active treatment ( $\beta_{W \cdot PI}$ ) by 2 (Supplementary Figure 1). The third “perturbed model” overestimates treatment effect heterogeneity by multiplying the coefficient of the prognostic index of control treatment ( $\beta_{(1-W) \cdot PI}$ ) by 2 and the coefficient of the prognostic index of active treatment ( $\beta_{W \cdot PI}$ ) by 0.5 (Supplementary Figure 1). We calculated the root mean squared error (RMSE) to indicate the level of perturbation for each model.

Finally, we computed the performance metrics for the “optimal” and the three “perturbed models” in the super population. We also visualized the performance of each of the four models with treatment effect calibration plots.

## 2.5 Case study

The performance of three different modelling approaches to predict treatment effect for patients at risk of diabetes in the DPP data set was compared using the proposed metrics.

The first approach (“risk model”) uses logistic regression to explain the outcome probability  $p_i = P(Y_i = 1 | X_i = x, W_i = w)$  based on the treatment indicator  $W$ , the centered prognostic index  $PI$  as defined before, and their interaction:

$$\log \frac{p_i}{1 - p_i} = W \cdot \beta_W + s(PI) \cdot \beta_{PI} + W \cdot s(PI) \cdot \beta_{W \cdot PI},$$

where  $s(\cdot)$  represents restricted cubic splines with two degrees of freedom.

The second approach (“effect model”) uses a penalized Ridge logistic regression to explain the outcome probability  $p_i$  based on the unpenalized treatment indicator  $W$ , penalized patient characteristics  $X$ , and their interaction:

$$\log \frac{p_i}{1 - p_i} = W \cdot \beta_W + X \cdot \beta_X + W \cdot X \cdot \beta_{W \cdot X}$$

where the level of penalization was determined by the minimum squared error in 5-fold cross-validation[18].

The third approach is a causal forest, which is similar to a random forest but maximizes heterogeneity in treatment effect rather than variation in the outcome[19]. Causal trees were built honestly by partitioning the data into two subsamples. One subsample was used to construct the trees, and another subsample to predict the treatment effect[19]. We used 1000 trees to tune the parameters (e.g., minimal node size, pruning) and 2000 trees to construct the final causal forest.

The models were trained on 70 percent of the patient data. The remaining 30 percent of the patient data, the test set, was used to calculate performance metrics with confidence intervals using 100 bootstrap samples of matched patient pairs. We used the R packages MatchIt for matching patients, mice for single imputation, stats for local regression, rms for restricted cubic splines, glmnet for Ridge penalization, and grf for causal forest (R version 4.1.0)[17, 20-24].

### **3. Results**

#### **3.1 Patient data**

Between 1996 and 2001, the DPP collected data on 3,081 participants of which 1,024 received lifestyle intervention, 1,027 received metformin, and 1,030 received placebo treatment (Supplementary Table 1). The median age of the participants was 52 years (IQR: 42-57 years), 67% of the participants were female, and the median

BMI value was 33 (IQR: 29-37). The proportion of patients developing diabetes was 4.8%, 7.0%, and 9.5% among participants receiving lifestyle intervention, metformin, and placebo treatment, respectively (Supplementary Table 1).

### 3.2 Simulation study

As expected, the treatment effect predictions of the “optimal model” were almost perfectly calibrated (calibration-in-the-large=-0.001,  $E_{\text{avg-for-benefit}}=0.001$ ,  $E_{50\text{-for-benefit}}=0.001$ ,  $E_{90\text{-for-benefit}}=0.003$ , Figure 2A). The “optimal model” was well able to discriminate (C-for-benefit=0.655, Figure 2A) between patients with small treatment harm (ATE=-0.023 in the quantile of patients with smallest predicted treatment effect) and patients with substantial treatment benefit (ATE=0.385 in the quantile of patients with largest treatment effect).

The first “perturbed model” was designed to overestimate treatment effect of lifestyle intervention (RMSE=0.090), which was expressed graphically by the corresponding calibration curve lying below the 45-degree line, and numerically by suboptimal calibration metrics (calibration-in-the-large=-0.074,  $E_{\text{avg-for-benefit}}=0.074$ ,  $E_{50\text{-for-benefit}}=0.065$ ,  $E_{90\text{-for-benefit}}=0.115$ , Figure 2B). The C-for-benefit expressed a slightly poorer ability to distinguish between patients with small and large treatment effects than the “optimal model” (C-for-benefit=0.649 versus 0.655). The cross-entropy-for-benefit and Brier-for-benefit also expressed poorer overall performance than the “optimal model” (cross-entropy-for-benefit=0.757 versus 0.733, Brier-for-benefit=0.215 versus 0.212, Figure 2A; 2B).

The second “perturbed model” was designed to overestimate risk heterogeneity of patients receiving lifestyle intervention (RMSE=0.058), which was expressed graphically by the corresponding calibration curve lying above the

diagonal for low predicted treatment effect (underestimation of low treatment effect) and below the diagonal for high predicted treatment effect (overestimation of high treatment effect), and numerically by suboptimal calibration metrics (calibration-in-the-large=-0.004,  $E_{\text{avg-for-benefit}}=0.070$ ,  $E_{50\text{-for-benefit}}=0.041$ ,  $E_{90\text{-for-benefit}}=0.189$ , Figure 2C). The C-for-benefit expressed a slightly poorer ability to distinguish between patients with small and large treatment effects than the “optimal model” (C-for-benefit=0.650 versus 0.655). The cross-entropy-for-benefit and Brier-for-benefit also expressed poorer overall performance than the “optimal model” (cross-entropy-for-benefit=0.784 versus 0.733, Brier-for-benefit=0.224 versus 0.212, Figure 2A; 2C).

The third “perturbed model” was designed to overestimate treatment effect heterogeneity of patients receiving lifestyle intervention (RMSE=0.164), which was expressed graphically by the corresponding calibration curve lying more extremely above the diagonal for low predicted treatment effect (underestimation of low treatment effect) and more extremely below the diagonal for high predicted treatment effect (overestimation of high treatment effect), and numerically by suboptimal calibration metrics ( $E_{\text{avg-for-benefit}}=0.124$ ,  $E_{50\text{-for-benefit}}=0.117$ ,  $E_{90\text{-for-benefit}}=0.230$ , Figure 2D). The C-for-benefit expressed a slightly poorer ability to distinguish between patients with small and large treatment effects than the “optimal model” (C-for-benefit=0.642 versus 0.655, Figure 2D). The cross-entropy-for-benefit and Brier-for-benefit also expressed poorer overall performance than the “optimal model” (cross-entropy-for-benefit=0.787 versus 0.733, Brier-for-benefit=0.221 versus 0.212, Figure 2A; 2D).

The results from the simulations using the metformin treatment arm rather than the lifestyle intervention arm were similar (Supplementary Figure 2; 3).

### 3.3 Case study

The differences in any of the performance measures between the risk model, the effect model, and the causal forest were not significantly different from zero in the 30 percent of patients who were in the test dataset ( $n=617$ ; Supplementary Table 1). Numerically, most calibration metrics of the effect model were better than that of the risk model (calibration-in-the-large=0.043 versus 0.052;  $E_{\text{avg-for-benefit}}=0.050$  versus 0.053;  $E_{90\text{-for-benefit}}=0.123$  versus 0.141, Figure 3A; 3B). Consequently, the overall performance of the effect model was numerically better than that of the risk model (cross-entropy-for-benefit=0.743 versus 0.747, Figure 3A; 3B), despite the numerically poorer discriminative ability of the effect model ( $C\text{-for-benefit}=0.663$  versus 0.664, Figure 3A; 3B).

Central calibration metrics of the causal forest were numerically poorer than those of the risk model ( $E_{\text{avg-for-benefit}}=0.074$  versus 0.053;  $E_{50\text{-for-benefit}}=0.068$  versus 0.031, Figure 3A; 3C), but the causal forest resulted in less extreme miscalibration than the risk model ( $E_{90\text{-for-benefit}}=0.101$  versus 0.141, Figure 3A; 3C). Due to less extreme miscalibration and numerically better discriminative ability ( $C\text{-for-benefit}=0.677$  versus 0.664, Figure 3A; 3C), the overall performance of the causal forest was numerically better than that of the risk model (cross-entropy-for-benefit=0.738 versus 0.747, Figure 3A; 3C).

## 4. Discussion

We extended the E-statistics, cross-entropy, and Brier score to quantify the quality of treatment effect predictions. The simulation study showed that the proposed metrics may be useful for comparing models because the metrics of the data-generating model were consistently better than those of deliberately “perturbed models”. The

case study illustrated the use of the proposed metrics in practice and showed that the calibration, discriminative ability, and overall performance of the three different models predicting treatment effect were not significantly different.

Similar to the previously proposed C-for-benefit, we defined observed pairwise treatment effect by the difference between outcomes in pairs of matched patients[6]. However, matching patients based on predicted treatment effect would result in different patient pairs and consequently different observed pairwise treatment effect for each prediction model[6]. Therefore, we chose to match patients based on the Mahalanobis distance between patient characteristics resulting in the same observed pairwise treatment effect for each prediction model. Alternative matching procedures, even matching patients randomly, resulted in similar conclusions in the simulation study (Supplementary Figure 6; 7). We matched without replacement since the treatment arms were similar in size, but matching with replacement is more appropriate for samples with unbalanced treatment arms. Furthermore, we selected relevant patient characteristics based on clinical expertise and existing literature, but variable selection is more suitable in high-dimensionality data.

The case study is merely an illustration of the use of the performance metrics and not a framework for model selection or internal validation. The use of internal validation techniques other than split sampling is recommended for quantification of the performance of a model in similar settings, but that was outside the scope of this study[25]. The choice of the percentage of observations used for the training and test set was arbitrary. Furthermore, the proposed metrics in the training set will not be insightful when using models with penalization and honest tree building, because they will indicate by definition miscalibration in the training set (Supplementary Figure 4; 5).

The strength of our study is that we propose currently lacking performance metrics for models predicting treatment effect. Their actual values can be used to compare models predicting treatment effect. Furthermore, in future research updating strategies can be considered if our proposed calibration metrics indicate miscalibration of treatment effect predictions.

A limitation of this study is the limited sample size of the case study. In the simulation study, we showed that the performance metrics were able to distinguish between models for an artificially enlarged data set. However, in the case study, the confidence intervals of the performance metrics were overlapping. This phenomenon is inherent to treatment effect estimation. To obtain reasonable power, heterogeneous treatment effect analyses require a much larger sample size compared to when estimating an overall ATE[26]. The case study suggested that there is a trade-off between calibration and discrimination: better calibrated models were worse at discriminating between patients with small and large treatment effects, but due to the small sample size no strict conclusions can be drawn. Secondly, the performance metrics were developed for binary outcomes, which could be extended to continuous outcomes in future research. Notwithstanding these limitations, we conclude that the proposed metrics are useful to assess the calibration and overall performance of models predicting treatment effect.



## **Author contributions**

**C.C.H.M. Maas:** Conceptualization, Methodology, Software, Formal analysis, Writing – Original Draft  
**D.M. Kent:** Conceptualization, Methodology, Writing – Review & Editing  
**M.C. Hughes:** Conceptualization, Writing – Review & Editing  
**R. Dekker:** Conceptualization, Writing – Review & Editing  
**H.F. Lingsma:** Conceptualization, Writing – Review & Editing  
**D. van Klaveren:** Conceptualization, Methodology, Writing – Review & Editing

## **Competing interests**

All authors declare no competing interests.

## **Data Availability**

Information on the process of obtaining the study dataset is available at the NIDDK Repository website (<https://repository.niddk.nih.gov/studies/dpp/>). The dataset can be obtained by submitting a formal request to the NIDDK Repository. The code to compute the results is available on <https://github.com/CHMMAas/PaperPredictionMetrics>.

## **Funding statement**

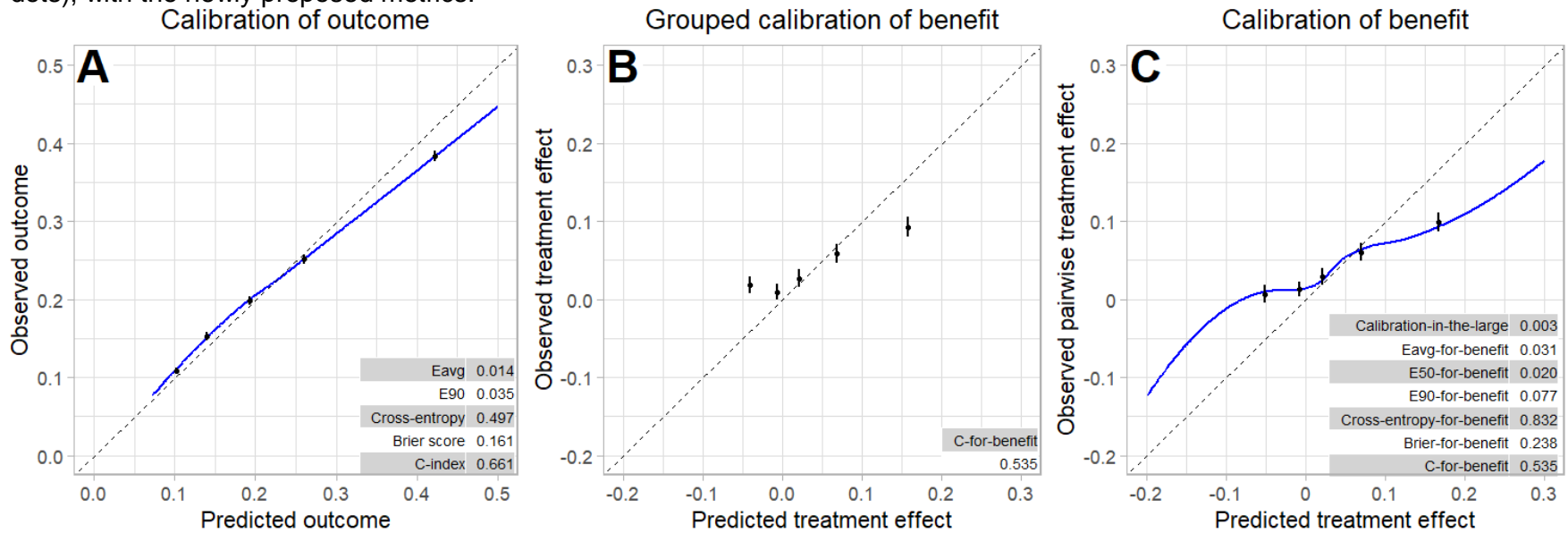
This work was partially supported through a Patient-Centered Outcomes Research Institute (PCORI) Award: the Predictive Analytics Resource Center (PARC) [SA.Tufts.PARC.OSCO.2018.01.25].

## References

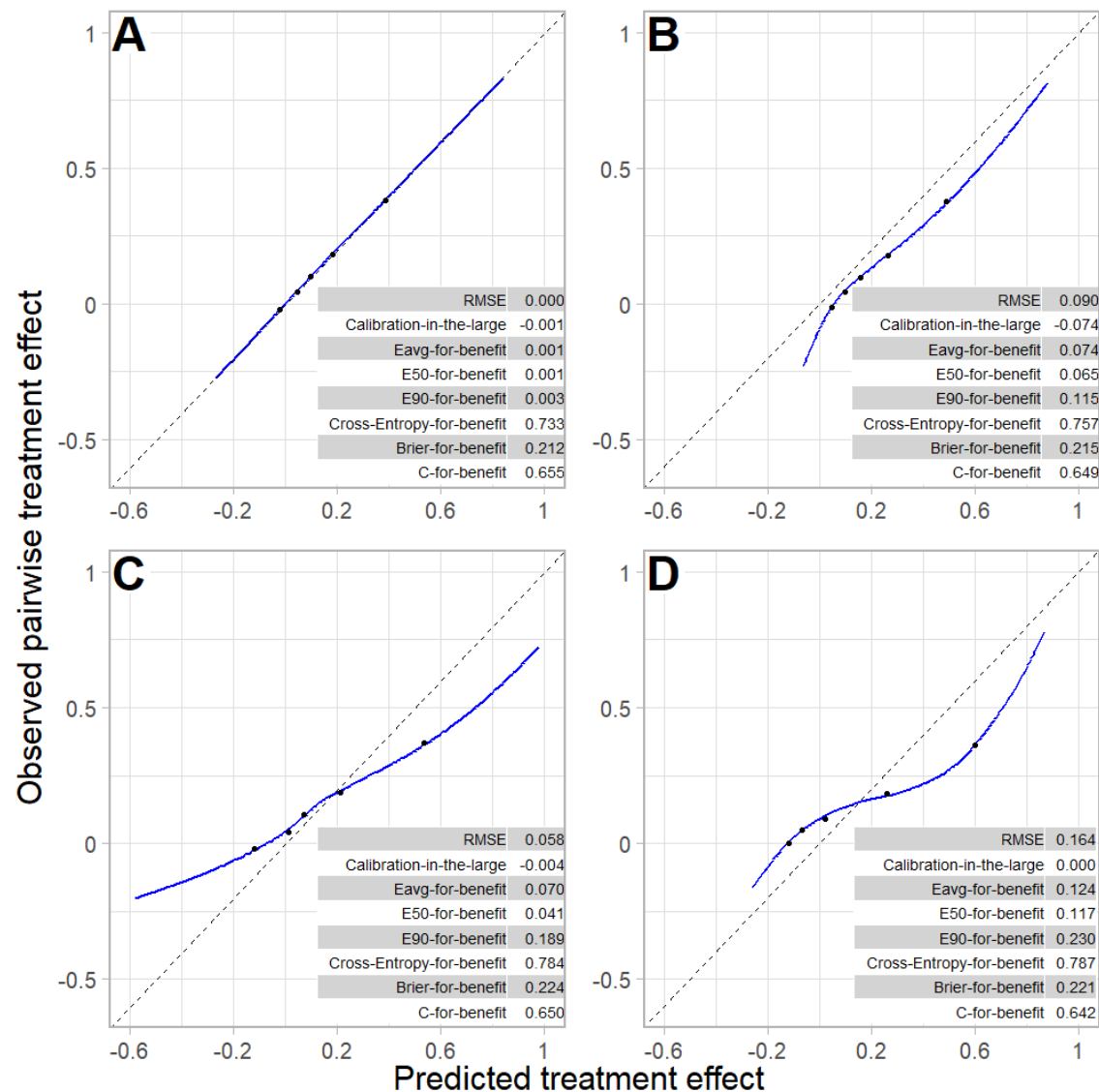
- [1] Rothwell PM. Can overall results of clinical trials be applied to all patients? *The Lancet*. 1995;345:1616-19.
- [2] Kravitz RL, Duan N, Braslow J. Evidence-Based Medicine, Heterogeneity of Treatment Effects, and the Trouble with Averages. *The Milbank Quarterly*. 2004;82:611-87.
- [3] Kent DM, Steyerberg EW, van Klaveren D. Personalized evidence based medicine: predictive approaches to heterogeneous treatment effects. *The BMJ*. 2018;363.
- [4] Rekkas A, Paulus JK, Raman G, Wong JB, Steyerberg EW, Rijnbeek PR, et al. Predictive approaches to heterogeneous treatment effects: a scoping review. *BMC Med Res Methodol*. 2020;20:264.
- [5] Tajik P, Oude Rengerink K, Mol BW, Bossuyt PM. SYNTAX score II. *The Lancet*. 2013;381:1899.
- [6] van Klaveren D., Steyerberg E.W., Serruys P.W., Kent D.M. The proposed 'concordance-statistic for benefit' provided a useful metric when modeling heterogeneous treatment effects. *J Clin Epidemiol*. 2018;94:59-68.
- [7] Austin P.C., Steyerberg E.W. The Integrated Calibration Index (ICI) and related metrics for quantifying the calibration of logistic regression models. *Stat Med*. 2019;38:4051-65.
- [8] Good I.J. Some Terminology and Notation in Information Theory. *Proceedings of the IEE - Part C: Monographs*. 1956;103:200-4.
- [9] Brier GW. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*. 1950;78.
- [10] van Klaveren D, Balan TA, Steyerberg EW, Kent DM. Models with interactions overestimated heterogeneity of treatment effects and were prone to treatment mistargeting. *Journal of Clinical Epidemiology*. 2019;114:72-83.
- [11] Takahashi K, van Klaveren D, Steyerberg EW, Onuma Y, Serruys PW. Concerns with the new SYNTAX score – Authors' reply. *The Lancet*. 2021;397:795-6.
- [12] Rubin D. Estimating Causal Effects of Treatments in Randomized and Nonrandomized Studies. *Journal of Educational Psychology*. 1974;66:688-701.
- [13] Ho DE, Imai K, King G, Stuart EA. Matching as Nonparametric Preprocessing for Reducing Model Dependence in Parametric Causal Inference. *Political Analysis*. 2007;15:199-236.
- [14] Van Calster B, McLernon DJ, van Smeden M, Wynants L, Steyerberg EW, Topic Group 'Evaluating diagnostic t, et al. Calibration: the Achilles heel of predictive analytics. *BMC Med*. 2019;17:230.
- [15] Maas CCHM. HTEPredictionMetrics: Heterogeneous Treatment Effect Prediction Metrics. R package version 1.0 available at: <https://github.com/CHMMAas/HTEPredictionMetrics> ed. Available at: <https://github.com/CHMMAas/HTEPredictionMetrics2022>.
- [16] Sussman J.B., Kent D.M., Nelson J.P., Hayward R.A. Improving diabetes prevention with benefit based tailored treatment: risk based reanalysis of Diabetes Prevention Program. *BMJ*. 2015;350:h454.
- [17] van Buuren S., Groothuis-Oudshoorn K. mice: Multivariate Imputation by Chained Equations in R. Available at: <https://www.jstatsoft.org/v45/i03/>; *Journal of Statistical Software*; 2011.
- [18] van Klaveren D., Vergouwe Y., Farooq V., Serruys P.W., Steyerberg E.W. Estimates of absolute treatment benefit for individual patients required careful modeling of statistical interactions. *J Clin Epidemiol*. 2015;68:1366-74.
- [19] Athey S., Imbens G. Recursive partitioning for heterogeneous causal effects. *Proc Natl Acad Sci U S A*. 2016;113:7353-60.
- [20] Ho D.E., Imai K., King G., Stuart E.A. MatchIt: Nonparametric Preprocessing for Parametric Causal Inference. 2011;42:1-28.
- [21] Team RC. R: A Language and Environment for Statistical Computing. Available at: <https://www.R-project.org/>; R Foundation for Statistical Computing; 2021.
- [22] Harrell F. E. Jr. rms: Regression Modeling Strategies. R package version 6.0-0. Available at: <https://CRAN.R-project.org/package=rms2020>.

- [23] Friedman J., Hastie T., R. T. Regularization Paths for Generalized Linear Models via Coordinate Descent. 2010;33:1-22.
- [24] Tibshirani J., Athey S., S. W. grf: Generalized Random Forests. R package version 1.2.0. Available at <https://CRAN.R-project.org/package=grf2020>.
- [25] Steyerberg EW. Clinical Prediction Models: A practical Approach to Development, Validation, and Updating: New York: Springer; 2009.
- [26] Brookes ST, Whitely E, Egger M, Smith GD, Mulheran PA, Peters TJ. Subgroup analyses in randomized trials: risks of subgroup-specific analyses; power and sample size for the interaction test. J Clin Epidemiol. 2004;57:229-36.

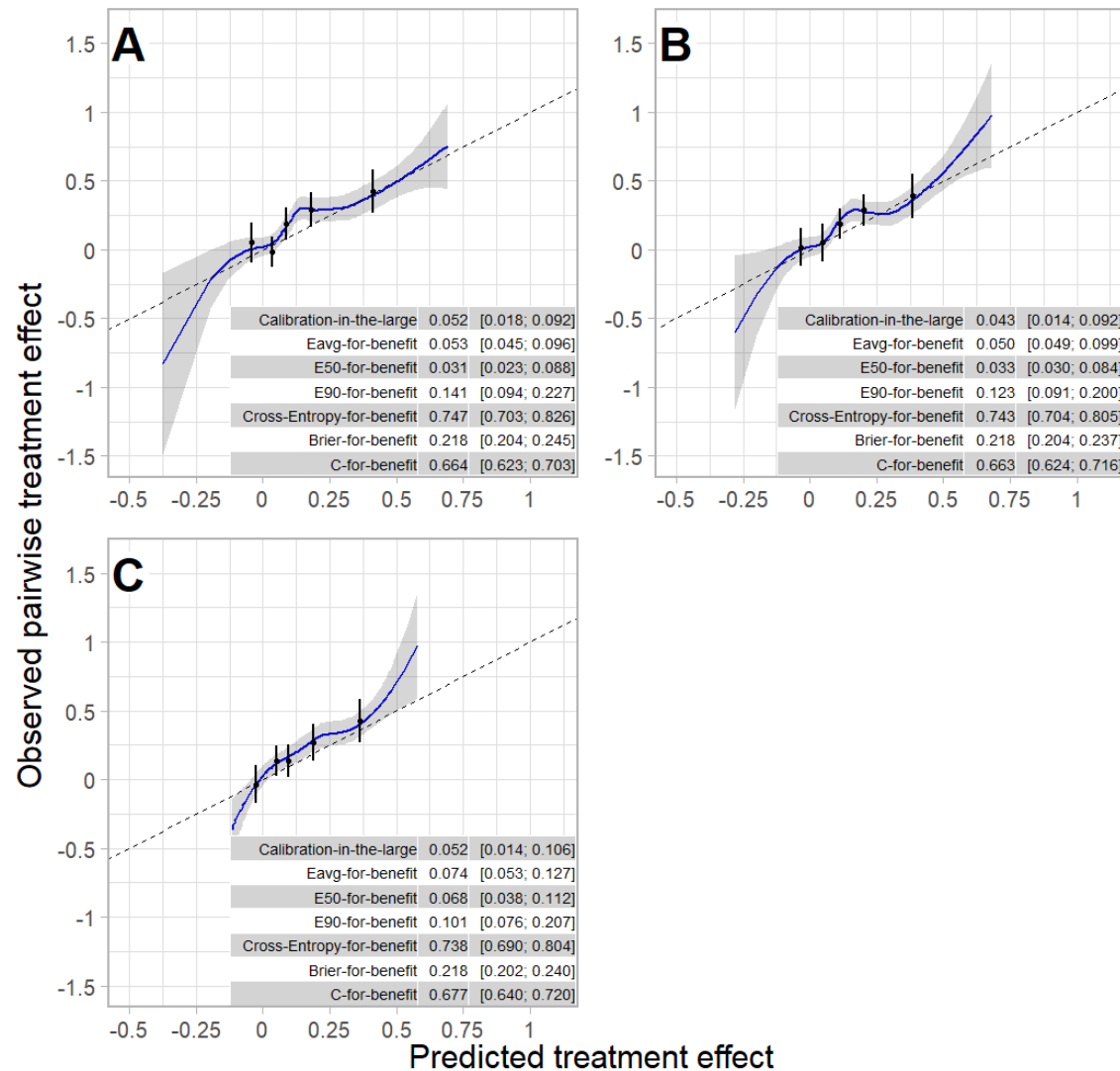
**Figure 1. Illustration of risk and benefit calibration figures with performance metrics of simulated data.** We sampled (n=3,600) from a simulated trial super population (1,000,000) with 12 binary risk predictors with 6 true treatment interactions[10]. Panel A depicts observed outcome versus predicted outcome by local regression (blue line, displayed between 0 and 0.5) and quantiles of predicted outcome (black dots), with the E-statistics, cross-entropy, Brier score, and C-index. Panel B depicts the calibration for benefit in groups with confidence intervals, with the C-for-benefit. Panel C depicts observed versus predicted treatment effect by local regression (blue line, displayed between -0.2 and 0.3) and quantiles of predicted treatment effect (black dots), with the newly proposed metrics.



**Figure 2. Calibration plot of the individualized treatment effect of simulated data from patients receiving lifestyle intervention.** This Figure depicts observed versus predicted treatment effect by smoothed calibration curves (blue line) and quantiles of predicted treatment effect (black dots) of simulated data from the lifestyle intervention versus placebo treatment. Observed pairwise treatment effect was obtained by matching patients based on patient characteristics. Smoothed calibration curves were obtained by local regression of the observed pairwise treatment effect of matched patient pairs on predicted treatment effect of matched patient pairs. For prediction of individualized treatment effect, we used a risk-based “optimal model” (panel **A**) and three “perturbed models” that overestimate average treatment effect (panel **B**), risk heterogeneity (panel **C**), and treatment effect heterogeneity (panel **D**). The average treatment effect is 12.9, 20.4, 12.9 (after a correction of -0.14), and 12.9 (after a correction of 0.53), respectively.



**Figure 3. Calibration plot of the individualized treatment effect of simulated data from patients receiving lifestyle intervention.** This Figure depicts observed versus predicted treatment effect by smoothed calibration curves (blue line with 95% confidence interval displayed by grey shaded area) and quantiles of predicted treatment effect (black dots) of lifestyle intervention versus placebo treatment. Observed pairwise treatment effect was obtained by matching patients based on patient characteristics. Smoothed calibration curves were obtained by local regression of the observed pairwise treatment effect of matched patient pairs on predicted treatment effect of matched patient pairs. For prediction of individualized treatment effect, we used: a risk modelling approach (panel **A**), a treatment effect modelling approach (panel **B**), and a causal forest (panel **C**). Confidence intervals around metric values were obtained using 100 bootstrap samples.



## **Performance metrics for models designed to predict treatment effect**

C.C.H.M. Maas<sup>1</sup>, D.M. Kent<sup>2</sup>, M.C. Hughes<sup>2</sup>, R. Dekker<sup>3</sup>, H.F. Lingsma<sup>1</sup>, D. van Klaveren<sup>1,2</sup>

<sup>1</sup>Department of Public Health, Erasmus University Medical Centre, Rotterdam, Netherlands

<sup>2</sup>Predictive Analytics and Comparative Effectiveness Centre, Institute for Clinical Research and Health Policy Studies, Tufts Medical Centre, Boston, USA

<sup>3</sup>Econometric Institute, Erasmus University Rotterdam, Rotterdam, Netherlands

*Correspondence:* C.C.H.M. Maas, Erasmus University Medical Centre, Doctor Molewaterplein 40, 3015 GD Rotterdam, Netherlands, [c.h.m.maas@erasmusmc.nl](mailto:c.h.m.maas@erasmusmc.nl)

## SUPPLEMENTARY INFORMATION

### Supplementary Information 1. Illustration of introduced metrics based on matched patient pairs.

Matched patient pair (A)	Patient assigned to treatment				Patient assigned to control treatment			
	$p_{0*}$ (B)	$p_{1**}$ (C)	Predicted benefit (D=B-C)	Observed outcome (E)	$p_{0*}$ (F)	$p_{1**}$ (G)	Predicted benefit (H=F-G)	Observed outcome (I)
1	0.136	0.283	-0.147	1	0.162	0.307	-0.145	1
2	0.246	0.343	-0.097	0	0.218	0.319	-0.101	1
3	0.156	0.219	-0.063	1	0.142	0.203	-0.061	0
4	0.081	0.083	0.043	0	0.098	0.062	0.036	0
5	0.345	0.212	0.133	1	0.299	0.171	0.128	0
6	0.421	0.390	0.306	1	0.561	0.255	0.306	1

\*  $p_0 = P(Y = 1|W = 0)$ ; \*\*  $p_1 = P(Y = 1|W = 1)$ ;

Matched patient pair (A)	$p_0$ (J=F)	$p_1$ (K=C)	Predicted benefit (L=J-K)	Observed benefit (M=E-I)	LOESS curve (N)*
1	0.162	0.283	-0.121	0	0
2	0.218	0.343	-0.125	-1	-1
3	0.142	0.219	-0.077	1	1
4	0.098	0.083	0.015	0	0
5	0.299	0.212	0.087	1	1
6	0.561	0.390	0.171	0	0

\*  $N = \text{predict}(\text{loess}(M \sim L))$ , which results in the same values as the observed benefit (M), when rounded to three decimals, due to a small number of observations.

The calibration metrics:

Overall-calibration-for-benefit =  $\text{abs}(\text{mean}(M) - \text{mean}(N)) = 0$

$E_{\text{avg}}\text{-for-benefit} = \text{mean}(\text{abs}(L-N)) \approx 0.529$

$E_{50}\text{-for-benefit} = \text{median}(\text{abs}(L-N)) \approx 0.523$

$E_{90}\text{-for-benefit} = \text{quantile}(\text{abs}(L-N), 0.9) \approx 0.995$

The overall performance:

Cross-entropy-for-benefit =  $-\frac{1}{n_p} [I(M = 1) \cdot \log[(1 - K)J] + I(M = 0) \log[(1 - K)(1 - J) + K \cdot J] + I(M = -1) \log[K(1 - J)]] \approx 1.049$

Brier-for-benefit =  $\frac{1}{2n_p} [[(1 - K)J - I(M = 1)]^2 + [(1 - K)(1 - J) + K \cdot J - I(M = 0)]^2 + [K(1 - J) - I(M = -1)]^2] \approx 0.321$



## Supplementary Information 2. Derivation of proper scoring rules in treatment effect estimation.

### Derivation of the Brier-for-benefit

The Brier-for-benefit is defined as

$$\text{Brier-for-benefit} = \frac{1}{2n_p} \sum_{i=1}^{n_p} \sum_{c \in \{-1, 0, 1\}} (P(\tau_i = c) - I(\tau_i = c))^2,$$

where  $n_p$  indicates the number of pairs,  $\tau_i$  indicates the observed pairwise treatment effect in a matched pair  $i$ ,  $I(\tau_i = c)$  is an indicator function returning one when the observed pairwise treatment effect of matched pair  $i$  ( $\tau_i$ ) is equal to class  $c$ , and  $P(\tau_i = c)$  indicates the probability that the observed pairwise treatment effect of matched pair  $i$  is equal to class  $c$ . The Brier score is divided by two to ensure that it lies between zero and one because in the worst-case scenario you give the highest prediction (one) for the wrong class, which would give a Brier score of two.

Equivalently,

$$\text{Brier-for-benefit} = \frac{1}{2n_p} \sum_{i=1}^{n_p} \left[ (P(\tau_i = 1) - I(\tau_i = 1))^2 + (P(\tau_i = 0) - I(\tau_i = 0))^2 + (P(\tau_i = -1) - I(\tau_i = -1))^2 \right]$$

Since matched patient pairs are independent, it holds that

$$\begin{aligned} P(\tau_i = 1) &= P(Y_i(1) = 0, Y_i(0) = 1) \\ &= P(Y_i(1) = 0)P(Y_i(0) = 1) \\ &= (1 - p_{i,1})p_{i,0} \\ P(\tau_i = 0) &= P((Y_i(1) = 0, Y_i(0) = 0) \cap (Y_i(1) = 1, Y_i(0) = 1)) \\ &= P(Y_i(1) = 0)P(Y_i(0) = 0) + P(Y_i(1) = 1)P(Y_i(0) = 1) \\ &= (1 - p_{i,1})(1 - p_{i,0}) + p_{i,1}p_{i,0} \\ P(\tau_i = -1) &= P(Y_i(1) = 1, Y_i(0) = 0) \\ &= P(Y_i(1) = 1)P(Y_i(0) = 0) \\ &= p_{i,1}(1 - p_{i,0}), \end{aligned}$$

where  $Y_i(W_i) = \begin{cases} Y_i(0) & \text{if } W_i = 0 \\ Y_i(1) & \text{if } W_i = 1 \end{cases}$  with  $Y_i$  indicates the potential outcome for patient  $i$

and  $W_i$  indicates the binary indicator for treatment, and the outcome probabilities conditional on treatment

$$\begin{aligned} p_{i,1} &= P(Y_i(1) = 1) = P(Y_i = 1 | W_i = 1) \\ p_{i,0} &= P(Y_i(0) = 1) = P(Y_i = 1 | W_i = 0). \end{aligned}$$

As a result, the Brier-for-benefit can be expressed as

$$\begin{aligned} \text{Brier-for-benefit} &= \frac{1}{2n_p} \sum_{i=1}^{n_p} \left( (1 - p_{i,1})p_{i,0} - I(\tau_i = 1) \right)^2 \\ &+ \frac{1}{2n_p} \sum_{i=1}^{n_p} \left( (1 - p_{i,1})(1 - p_{i,0}) + p_{i,1}p_{i,0} - I(\tau_i = 0) \right)^2 \\ &+ \frac{1}{2n_p} \sum_{i=1}^{n_p} \left( p_{i,1}(1 - p_{i,0}) - I(\tau_i = -1) \right)^2. \end{aligned}$$

### *Derivation of the cross-entropy-for-benefit*

Similarly, the cross-entropy-for-benefit is defined as

$$\begin{aligned}
 \text{Cross-entropy-for-benefit} &= -\frac{1}{n_p} \cdot \sum_{i=1}^{n_p} \sum_{c \in \{-1, 0, 1\}} I(\tau_i = c) \log[P(\tau_i = c)] \\
 &= -\frac{1}{n_p} \cdot \sum_{i=1}^{n_p} I(\tau_i = 1) \log[(1 - p_{i,1})p_{i,0}] \\
 &\quad - \frac{1}{n_p} \cdot \sum_{i=1}^{n_p} I(\tau_i = 0) \log[(1 - p_{i,1})(1 - p_{i,0}) + p_{i,1}p_{i,0}] \\
 &\quad - \frac{1}{n_p} \cdot \sum_{i=1}^{n_p} I(\tau_i = -1) \log[p_{i,1}(1 - p_{i,0})].
 \end{aligned}$$

### *Outcome probabilities of the causal forest*

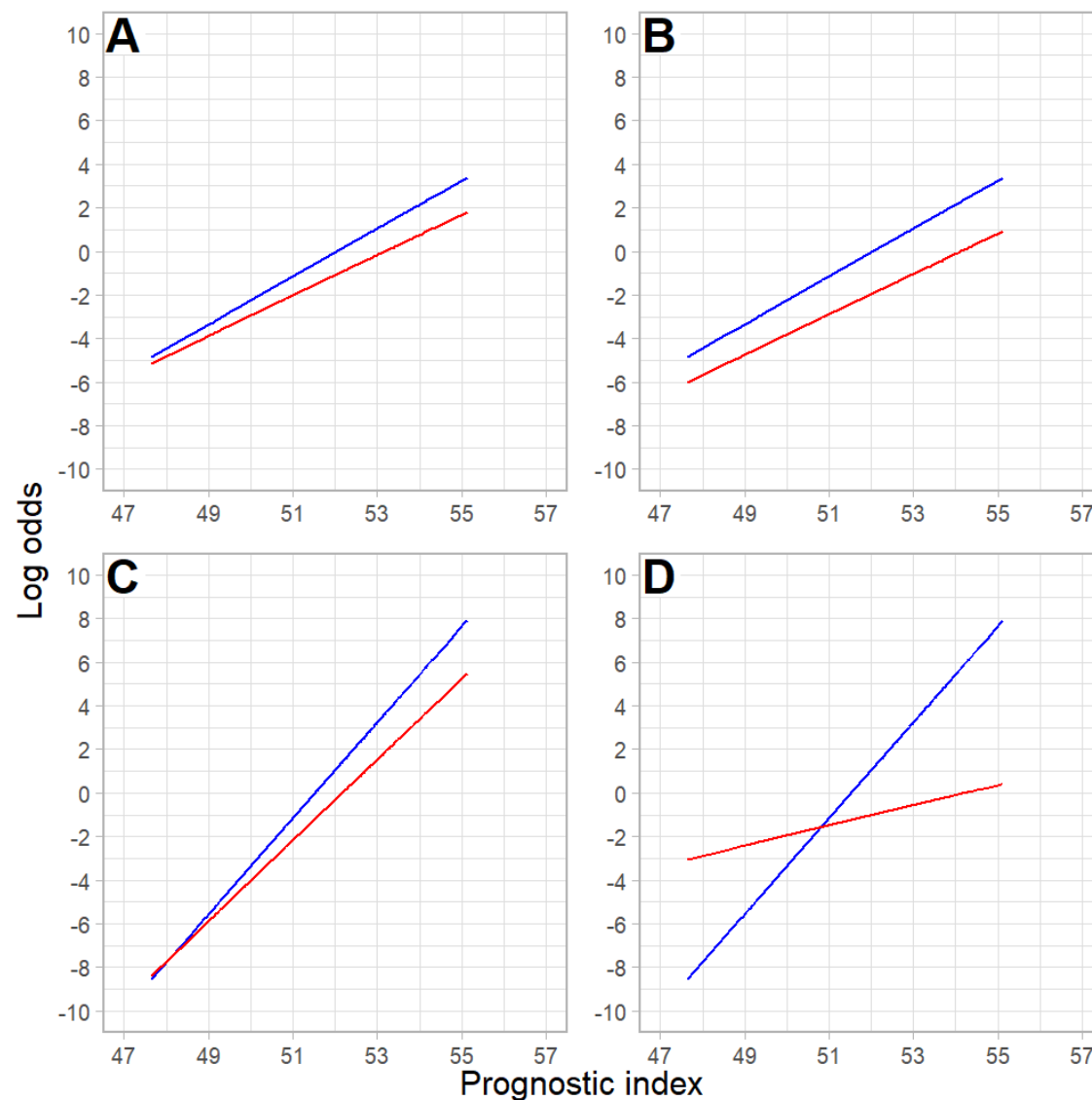
Of note, the outcome probabilities conditional on treatment  $p_{i,0}$  and  $p_{i,1}$  probabilities of the causal forest are obtained by

$$\begin{aligned}
 p_{i,0} &= E[Y|X, W = 0] = \hat{m}(X) - \hat{e}(X)\hat{\tau}(X) \\
 p_{i,1} &= E[Y|X, W = 1] = \hat{m}(X) + (1 - \hat{e}(X))\hat{\tau}(X),
 \end{aligned}$$

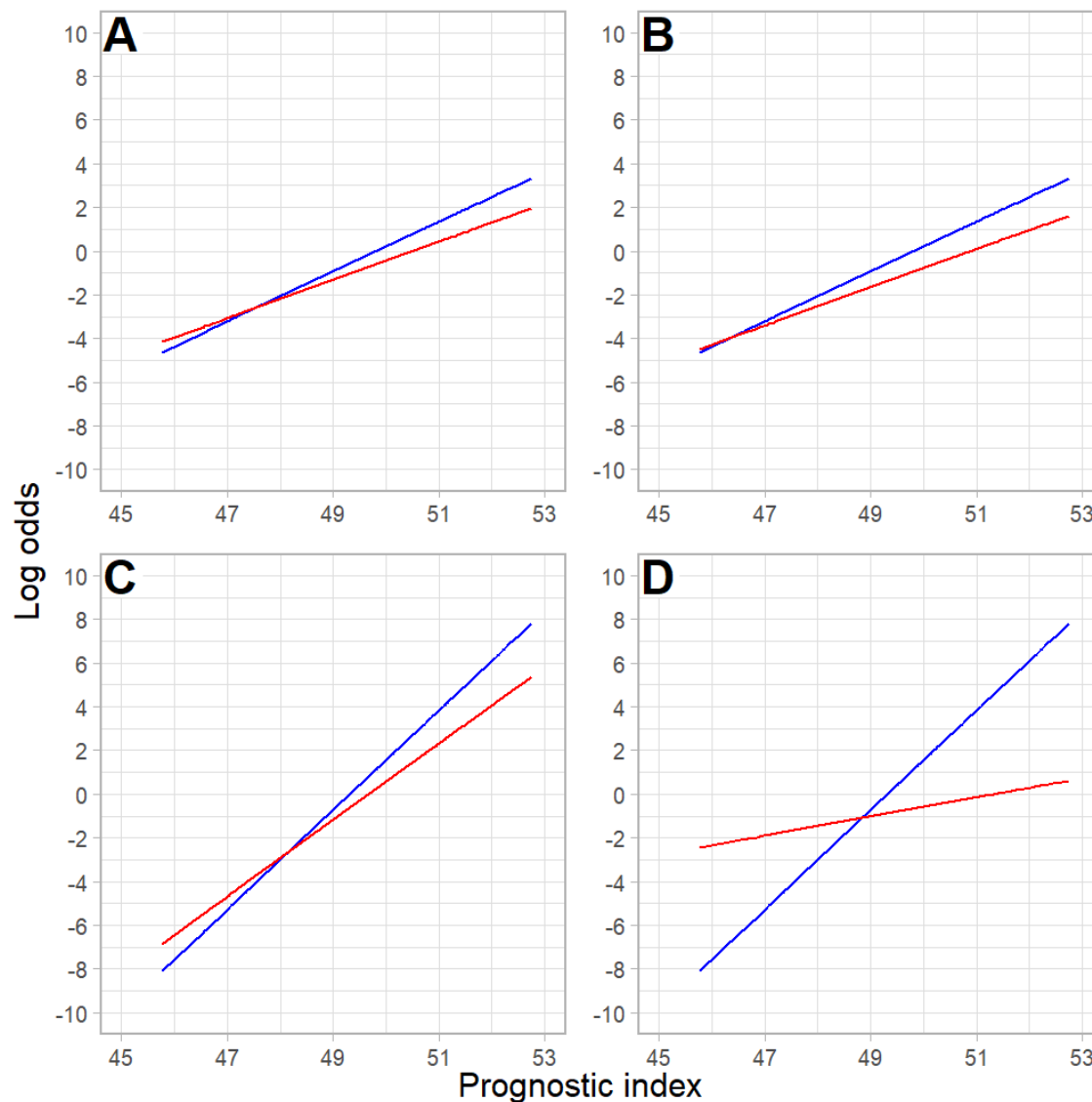
with  $\hat{m}(X) = E[Y|X = x]$  and  $\hat{e}(X) = E[W|X = x]$  indicate the outcomes of two random forests, and  $\hat{\tau}(X) = E[Y_i(0) - Y_i(1)|X = x]$  indicates the treatment effect outcomes.

## SUPPLEMENTARY FIGURES

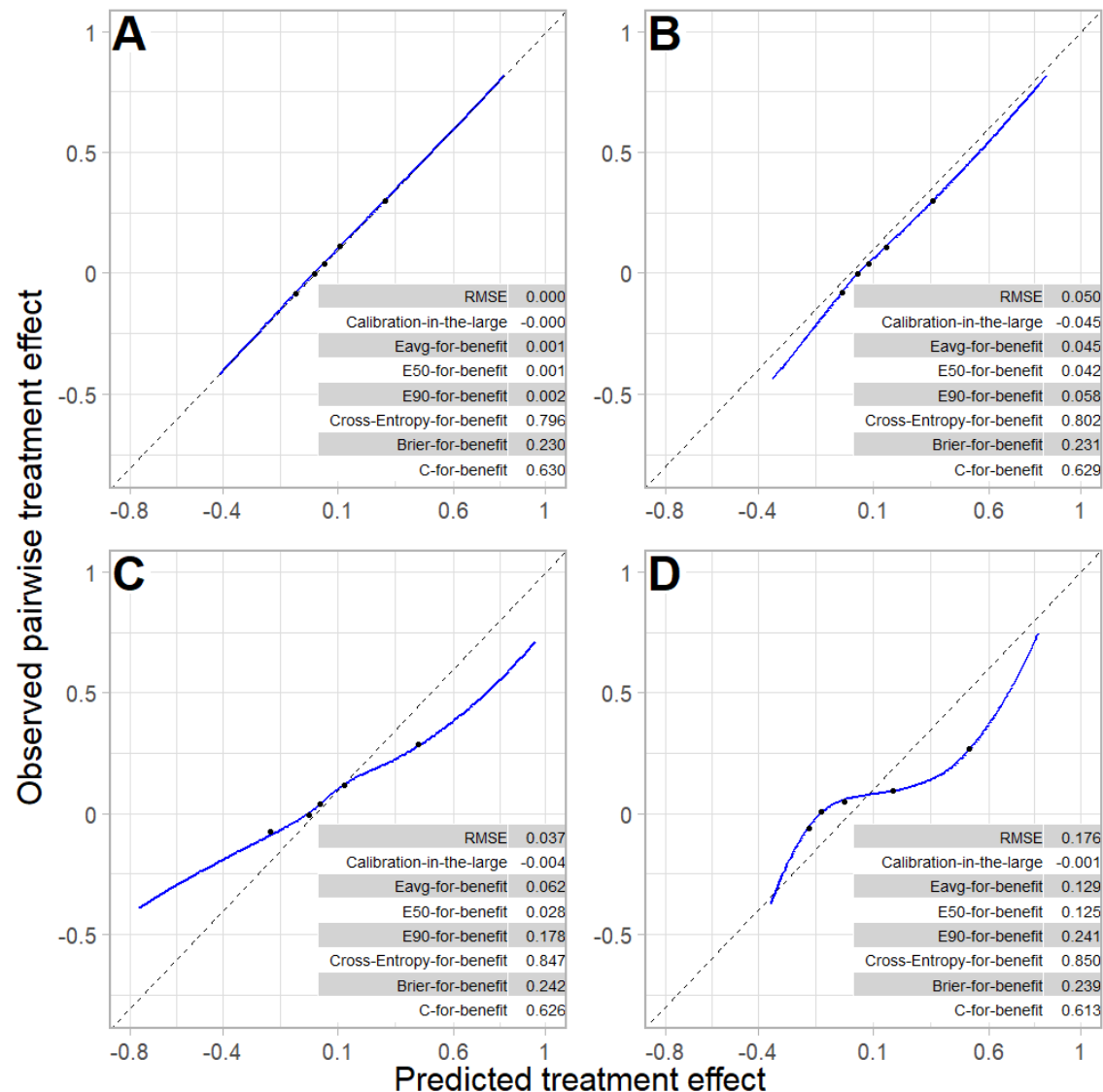
**Supplementary Figure 1. The probability of diabetes for each level of the prognostic index on the log odds scale when not treating (blue) and treating patients (red) with lifestyle intervention. This figure displays the “optimal model” in panel A, and three “perturbed models” that overestimate average treatment effect (panel B), risk heterogeneity (panel C), and treatment effect heterogeneity (panel D).**



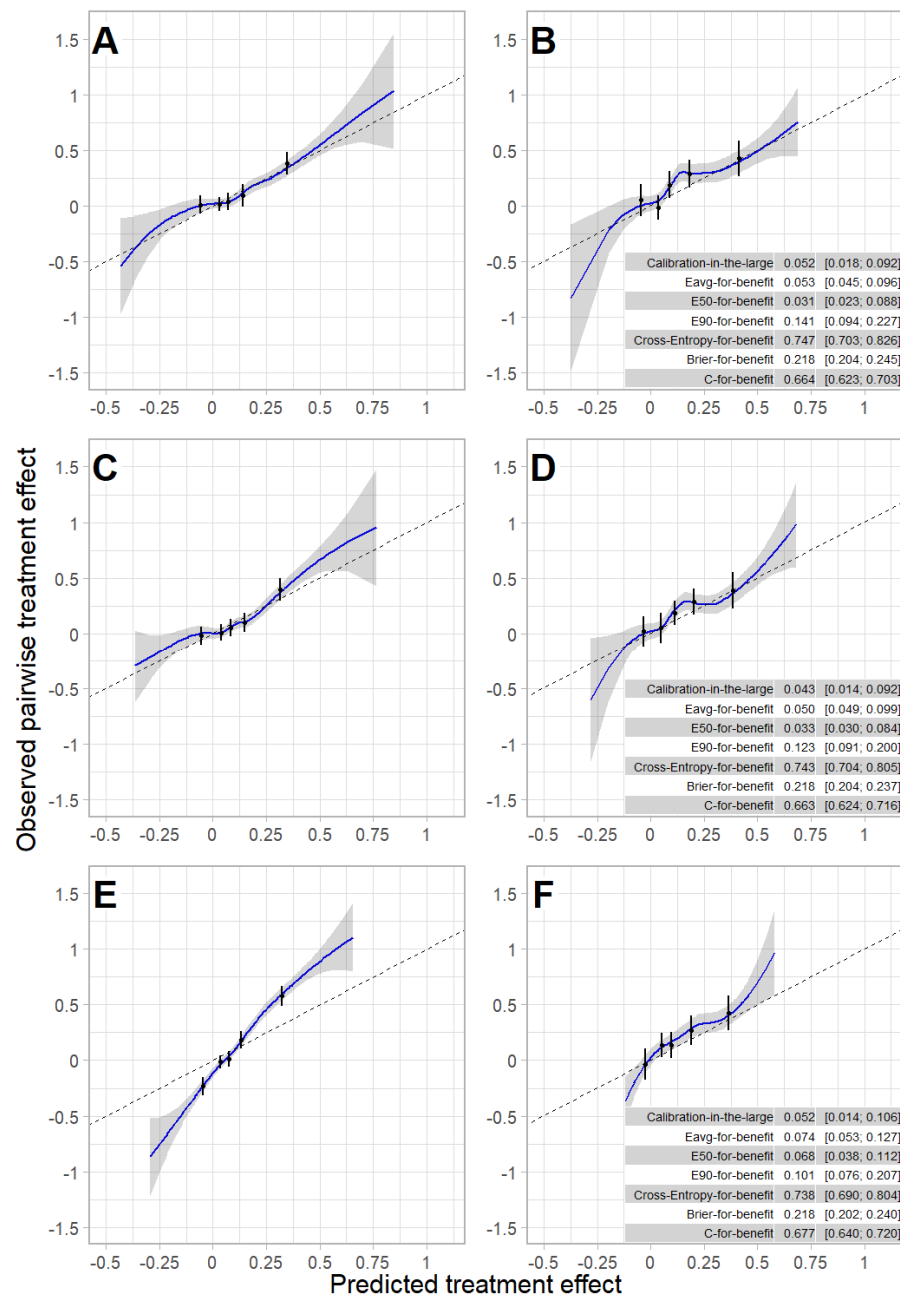
**Supplementary Figure 2. The probability of diabetes for each level of the prognostic index on the log odds scale when not treating (blue) and treating patients (red) with metformin. This figure displays the “optimal model” in panel A, and three “perturbed models” that overestimate average treatment effect (panel B), risk heterogeneity (panel C), and treatment effect heterogeneity (panel D).**



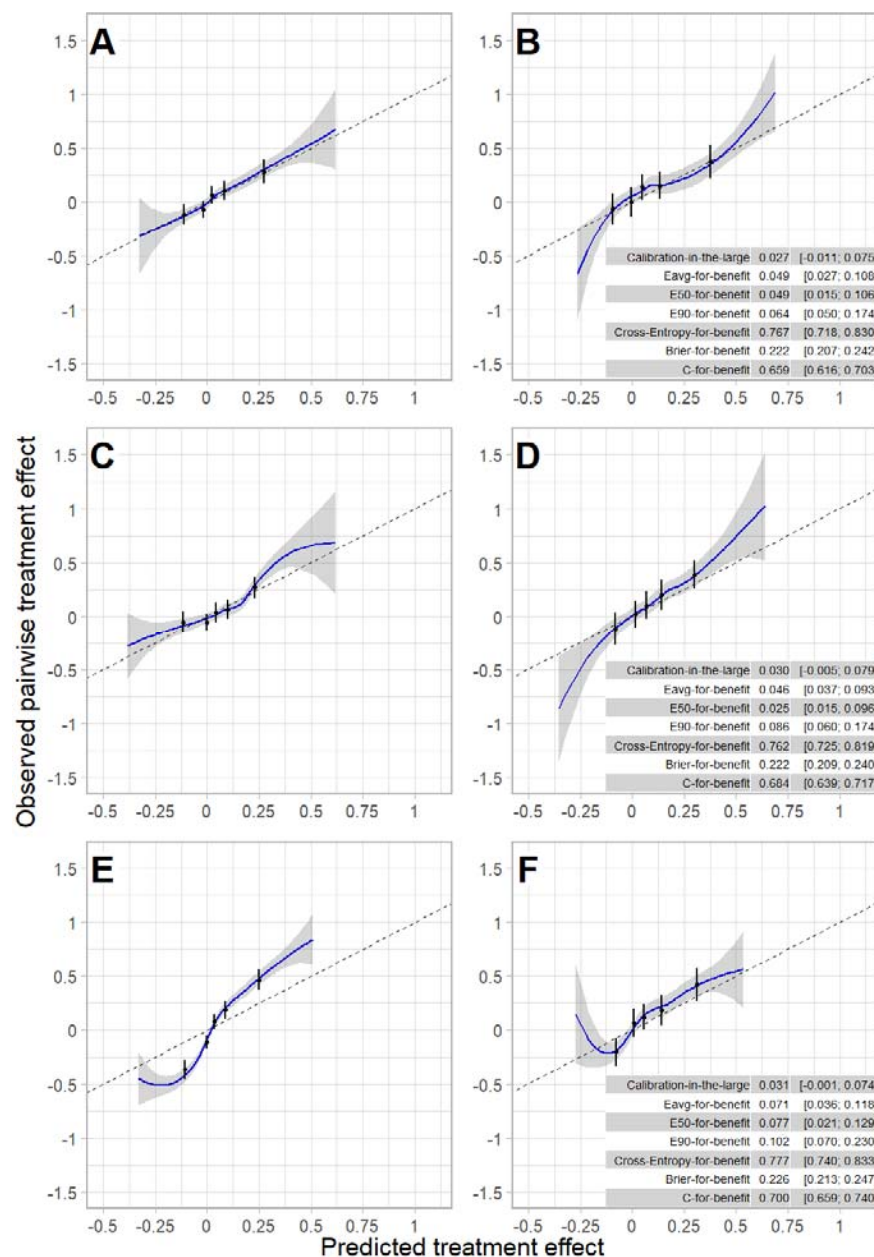
**Supplementary Figure 3. Calibration plot of the treatment effect of simulated data from patients receiving metformin.** This Figure depicts observed versus predicted treatment effect by smoothed calibration curves (blue line) and quarters of predicted treatment effect (black dots) of metformin versus placebo treatment. Observed pairwise treatment effect was obtained by matching patients based on patient characteristics. Smoothed calibration curves were obtained by local regression of the observed pairwise treatment effect of matched patient pairs on predicted treatment effect of matched patient pairs. For prediction of treatment effect, we used a risk-based “optimal model” (panel **A**) and three “perturbed models” that overestimate average treatment effect (panel **B**), risk heterogeneity (panel **C**), and treatment effect heterogeneity (panel **D**). The average treatment effect is 6.5, 11.1, 6.5 (after a correction of -0.02), and 6.5 (after a correction of 0.375), respectively.



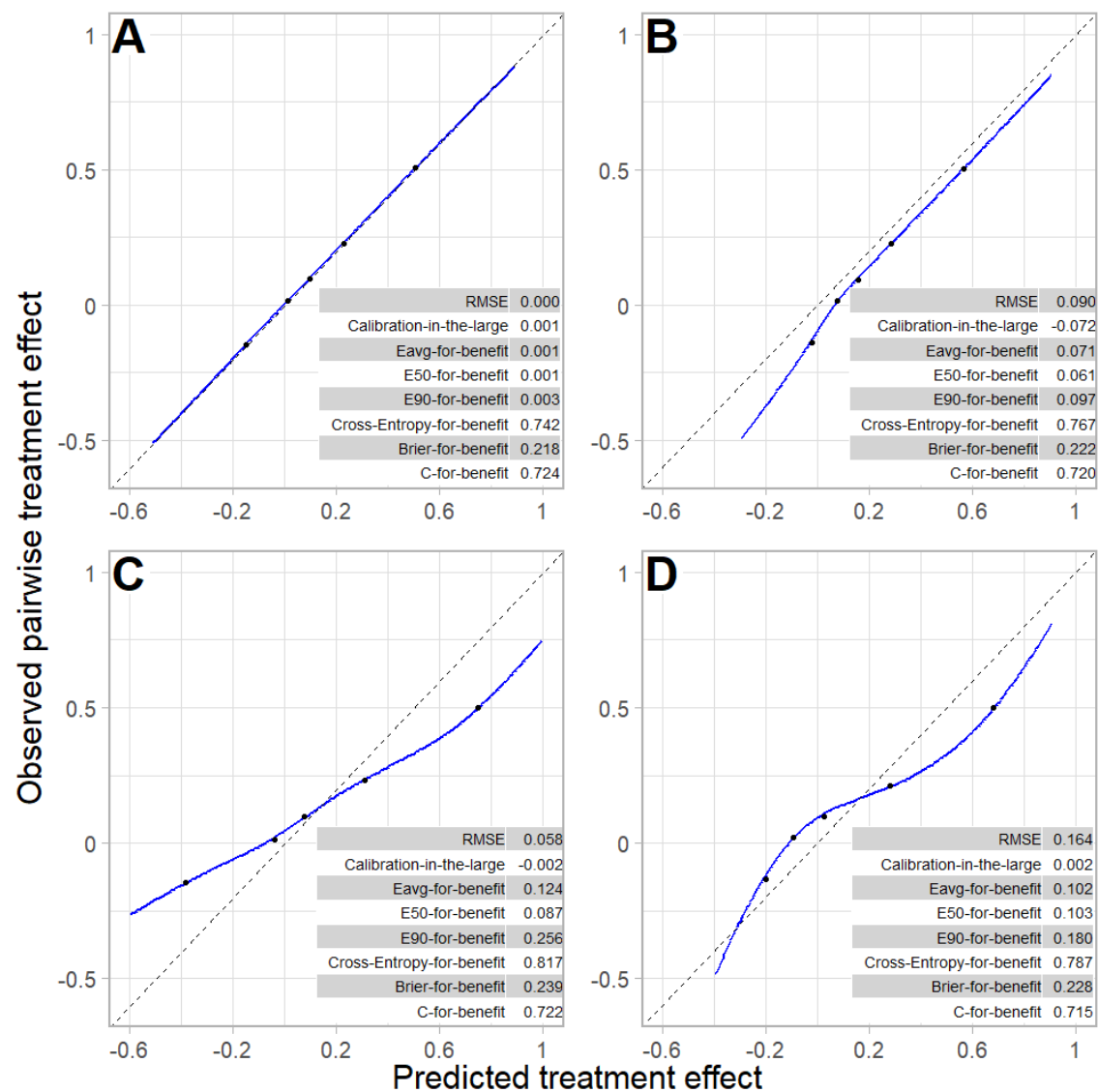
**Supplementary Figure 4. Calibration plot of the treatment effect of training and test data of lifestyle intervention.** This Figure depicts observed versus predicted treatment effect by smoothed calibration curves (blue with 95% confidence interval displayed by grey shaded area) and quarters of predicted treatment effect (black dots) of lifestyle intervention versus placebo treatment. Observed pairwise treatment effect was obtained by matching patients based on patient characteristics. Smoothed calibration curves were obtained by local regression of the observed pairwise treatment effect of matched patient pairs on predicted treatment effect of matched patient pairs. For prediction of treatment effect, we used: a risk modelling approach (panel **A**; **B**), a treatment effect modelling approach (panel **C**; **D**), and a causal forest (panel **E**; **F**). The models are trained on 70 percent of the data (panel **A**; **C**; **E**) and predictions are obtained on the other 30 percent of the data (**B**; **D**; **F**). Confidence intervals around the metric values were obtained using 100 bootstrap samples.



**Supplementary Figure 5. Calibration plot of the treatment effect of training and test data of metformin.** This Figure depicts observed versus predicted treatment effect by smoothed calibration curves (blue line with 95% confidence interval displayed by grey shaded area) and quarters of predicted treatment effect (black dots) of metformin versus placebo treatment. Observed pairwise treatment effect was obtained by matching patients based on patient characteristics. Smoothed calibration curves were obtained by local regression of the observed pairwise treatment effect of matched patient pairs on predicted treatment effect of matched patient pairs. For prediction of treatment effect, we used: a risk modelling approach (panel **A**; **B**), a treatment effect modelling approach (panel **C**; **D**), and a causal forest (panel **E**; **F**). The models are trained on 70 percent of the data (panel **A**; **C**; **E**) and predictions are obtained on the other 30 percent of the data (**B**; **D**; **F**). Confidence intervals around the metric values were obtained using 100 bootstrap samples.

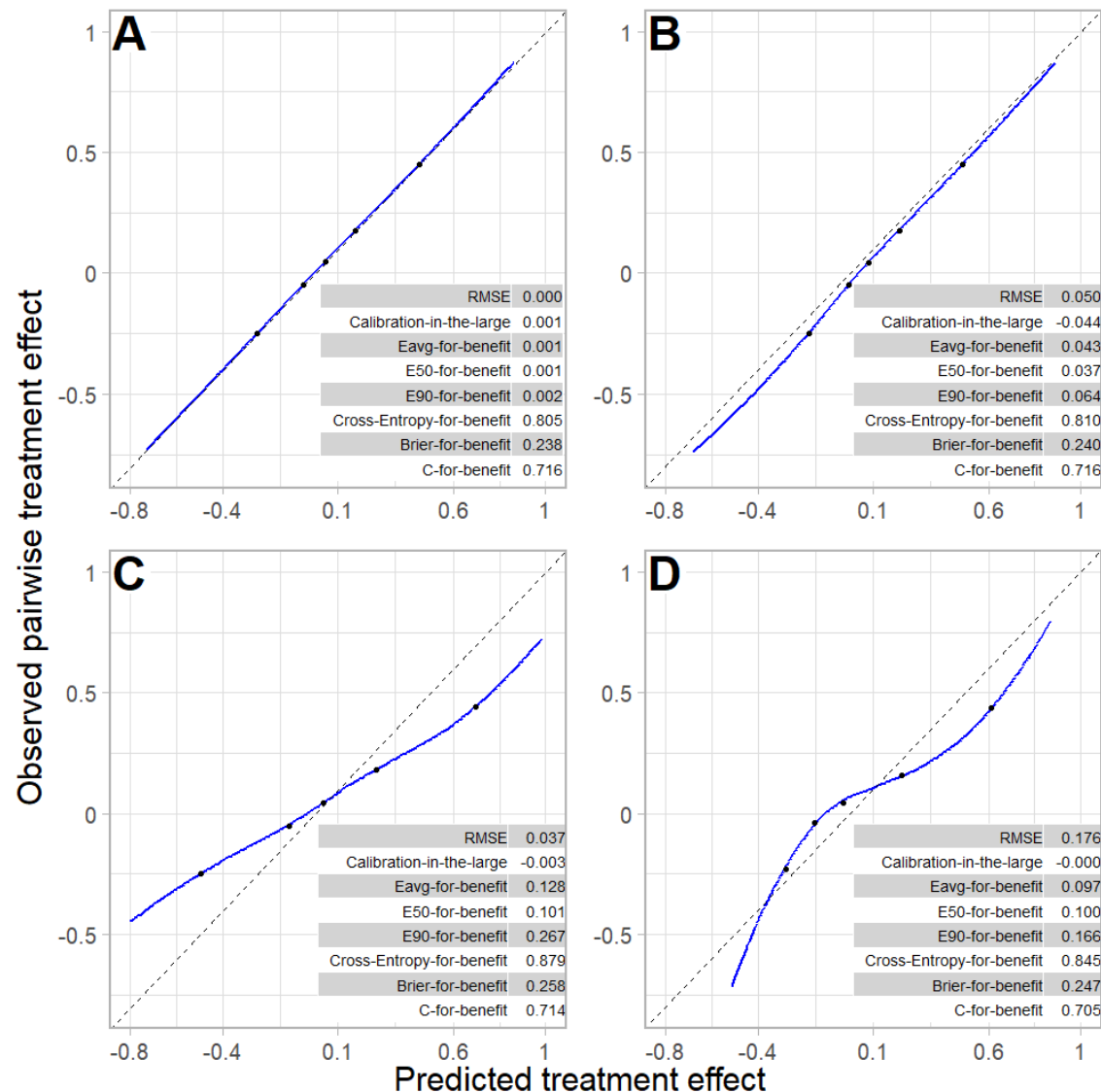


**Supplementary Figure 6. Calibration plot of the individualized treatment effect of simulated data from patients receiving lifestyle intervention.** This Figure depicts observed versus predicted treatment effect by smoothed calibration curves (blue line) and quantiles of predicted treatment effect (black dots) of simulated data from the lifestyle intervention versus placebo treatment. Observed pairwise treatment effect was obtained by matching patients *randomly*. Smoothed calibration curves were obtained by local regression of the observed pairwise treatment effect of matched patient pairs on predicted treatment effect of matched patient pairs. For prediction of individualized treatment effect, we used a risk-based “optimal model” (panel **A**) and three “perturbed models” that overestimate average treatment effect (panel **B**), risk heterogeneity (panel **C**), and treatment effect heterogeneity (panel **D**). The average treatment effect is 12.9, 20.4, 12.9 (after a correction of -0.14), and 12.9 (after a correction of 0.53), respectively.





**Supplementary Figure 7. Calibration plot of the treatment effect of simulated data from patients receiving metformin.** This Figure depicts observed versus predicted treatment effect by smoothed calibration curves (blue line) and quarters of predicted treatment effect (black dots) of metformin versus placebo treatment. Observed pairwise treatment effect was obtained by matching patients *randomly*. Smoothed calibration curves were obtained by local regression of the observed pairwise treatment effect of matched patient pairs on predicted treatment effect of matched patient pairs. For prediction of treatment effect, we used a risk-based “optimal model” (panel **A**) and three “perturbed models” that overestimate average treatment effect (panel **B**), risk heterogeneity (panel **C**), and treatment effect heterogeneity (panel **D**). The average treatment effect is 6.5, 11.1, 6.5 (after a correction of -0.02), and 6.5 (after a correction of 0.375), respectively



## SUPPLEMENTARY TABLES

**Supplementary Table 1. Characteristics of patients in the Diabetes Prevention Program receiving lifestyle intervention, metformin, and placebo treatment.**

	Total			Lifestyle		Metformin		Placebo	
	N	%	Missing	N	%	N	%	N	%
<b>Sample size</b>	3081			1024		1027		1030	
<b>Diabetes</b>	655	21.3		148	4.8	215	7	292	9.5
<b>Female</b>	2053	66.6		685	22.2	669	21.7	699	22.7
<b>Ethnicity</b>									
<b>Black</b>	644	20.9		204	6.6	221	7.2	219	7.1
<b>Hispanic</b>	508	16.5		178	5.8	162	5.3	168	5.5
<b>History of high blood glucose</b>	614	19.9		206	6.7	192	6.2	216	7
<b>Family history of diabetes</b>	2127	69	2	713	23.1	699	22.7	715	23.2
<b>Smoking</b>	216	7		67	2.2	69	2.2	80	2.6
<b>Hypertension</b>	835	27.1		286	9.3	267	8.7	282	9.2
<b>Gestational diabetes mellitus</b>	321	10.4	1	108	3.5	106	3.4	107	3.5
<b>Age</b>	52	[42; 57]		47	[42; 57]	52	[42; 57]	47	[42; 57]
<b>BMI</b>	33	[29; 37]		33	[29; 37]	33	[29; 37]	33	[29; 37]
<b>Triglycerides</b>	141	[99; 201]	5	138	[97; 200]	137	[98; 195]	147	[104; 208]
<b>Haemoglobin A<sub>1c</sub></b>	5.9	[5.6; 6.2]	8	5.9	[5.6; 6.2]	5.9	[5.6; 6.2]	5.9	[5.6; 6.2]
<b>Fasting blood sugar</b>	105	[101; 112]		105	[101; 112]	105	[100; 112]	106	[101; 112]

