

Online Supplement

Appendix A: Link of Prevented Diseases to Cost Minimization

Proof of Proposition 1. Without loss of generality, we study year l within our study period, and let $\{1, \dots, k\} \subset \{1, \dots, N\}$ denote the patients that were enrolled for preventive treatment according to our decision model, i.e., which describe a solution to the decision problem (Equation (6)). For simplicity, we neglect the index denoting the year. Then, the first part of this proof shows that the solution $\{1, \dots, k\}$ minimizes costs among all subsets containing k patients. The second part proves that there is no subset with fewer patients ($\{1, \dots, k-1\} \subset \{1, \dots, N\}$) that yields lower costs.

1st Part. Because $\{1, \dots, k\} \subset \{1, \dots, N\}$ is a solution to the decision problem, we can assume that

$$\begin{aligned} \mathbb{E}[y | x_i, t_i = 0] - \mathbb{E}[y | x_i, t_i = 1] &\geq \mathbb{E}[y | x_j, t_j = 0] - \mathbb{E}[y | x_j, t_j = 1], \\ \forall i \in \{1, \dots, k\} \forall j \in \{k+1, \dots, N\}. \end{aligned} \quad (12)$$

We take a random patient $i^* \in \{1, \dots, k\}$ and a random patient $j^* \in \{k+1, \dots, N\}$ and exchange the treatment (i.e., we provide preventive treatment to patient j^* but not to patient i^*). We compute the expected costs C_{after} after this exchange and compare it to the expected costs C_{before} before the exchange. The computation yields the following:

$$\begin{aligned} &C_{\text{after}} - C_{\text{before}} \quad (13) \\ &= C_{\text{diab}} \left(\sum_{\substack{i=1 \\ i \neq i^*}}^k \mathbb{E}[y | x_i, t_i = 1] + \sum_{\substack{j=k+1 \\ j \neq j^*}}^N \mathbb{E}[y | x_j, t_j = 0] + \mathbb{E}[y | x_{j^*}, t_{j^*} = 1] + \mathbb{E}[y | x_{i^*}, t_{i^*} = 0] \right) \\ &\quad + k C_{\text{prevent}} - C_{\text{diab}} \left(\sum_{i=1}^k \mathbb{E}[y | x_i, t_i = 1] + \sum_{j=k+1}^N \mathbb{E}[y | x_j, t_j = 0] \right) - k C_{\text{prevent}} \quad (14) \\ &= C_{\text{diab}} \left(\sum_{\substack{i=1 \\ i \neq i^*}}^k \mathbb{E}[y | x_i, t_i = 1] + \sum_{\substack{j=k+1 \\ j \neq j^*}}^N \mathbb{E}[y | x_j, t_j = 0] + \mathbb{E}[y | x_{j^*}, t_{j^*} = 1] + \mathbb{E}[y | x_{i^*}, t_{i^*} = 0] \right. \\ &\quad \left. - \sum_{i=1}^k \mathbb{E}[y | x_i, t_i = 1] - \sum_{j=k+1}^N \mathbb{E}[y | x_j, t_j = 0] \right) \quad (15) \\ &= C_{\text{diab}} \left(\mathbb{E}[y | x_{i^*}, t_{i^*} = 0] - \mathbb{E}[y | x_{i^*}, t_{i^*} = 1] + \mathbb{E}[y | x_{j^*}, t_{j^*} = 1] - \mathbb{E}[y | x_{j^*}, t_{j^*} = 0] \right) \quad (16) \\ &\geq 0 \quad (17) \end{aligned}$$

2nd Part. We now show that no subset comprising fewer patients ($\{1, \dots, k-1\} \subset \{1, \dots, N\}$) yields lower costs. For this computation, we assume that:

$$C_{\text{diab}} \mathbb{E}[y | x_i, t_i = 1] + C_{\text{prevent}} \leq C_{\text{diab}} \mathbb{E}[y | x_i, t_i = 0] \quad \forall i \in \{1, \dots, k\}. \quad (18)$$

That an exchange of patients between the treatment group and the non-treatment group leads to increased costs directly follows from part 1. Thus, without loss of generality, we consider subgroups $\{1, \dots, k-1\}$ of our original solution $\{1, \dots, k\}$ to the decision problem. We compute costs C_{k-1} when providing preventive care to $k-1$ patients and compare C_{k-1} to the original costs C_k :

$$C_{k-1} - C_k \quad (19)$$

$$= C_{\text{diab}} \left(\sum_{i=1}^{k-1} \mathbb{E}[y | x_i, t_i = 1] + \sum_{j=k}^N \mathbb{E}[y | x_j, t_j = 0] \right) + (k-1) C_{\text{prevent}}$$

$$- C_{\text{diab}} \left(\sum_{i=1}^k \mathbb{E}[y | x_i, t_i = 1] + \sum_{j=k+1}^N \mathbb{E}[y | x_j, t_j = 0] \right) - k C_{\text{prevent}} \quad (20)$$

$$= C_{\text{diab}} \left(\sum_{i=1}^{k-1} \mathbb{E}[y | x_i, t_i = 1] + \sum_{j=k}^N \mathbb{E}[y | x_j, t_j = 0] - \sum_{i=1}^k \mathbb{E}[y | x_i, t_i = 1] - \sum_{j=k+1}^N \mathbb{E}[y | x_j, t_j = 0] \right) - C_{\text{prevent}} \quad (21)$$

$$= C_{\text{diab}} \left(\mathbb{E}[y | x_k, t_k = 0] - \mathbb{E}[y | x_k, t_k = 1] \right) - C_{\text{prevent}} \quad (22)$$

$$= C_{\text{diab}} \left(\mathbb{E}[y | x_k, t_k = 0] \right) - C_{\text{diab}} \left(\mathbb{E}[y | x_k, t_k = 1] \right) - C_{\text{prevent}} \quad (23)$$

$$\geq 0 \quad (24)$$

□

Appendix B: Proof for Optimal Allocation

Proof of Proposition 2. Let n_{onsets} denote the number of expected onsets. Assume $h_0^*(x_{i,l})$ to be a perfect model for $\mathbb{E}[y | x_{i,l}, t_{i,l} = 0]$ and assume $\gamma_{i,l}$ to be the true treatment effect for $i \in \{1, \dots, N\}$. Let $\{1, \dots, k\}$ be the subset of patients who are prescribed preventive care, and let $\{k+1, \dots, N\}$ be the subset of patients who are not prescribed preventive care in year l . Let this allocation follow our described approach to allocating patients; thus,

$$(1 - \gamma_{i,l}) h_0^*(x_{i,l}) \geq (1 - \gamma_{j,l}) h_0^*(x_{j,l}), \quad \forall i \in \{1, \dots, k\} \forall j \in \{k+1, \dots, N\}, \quad (25)$$

which is equal to

$$h_0^*(x_{i,l}) + \gamma_{j,l} h_0^*(x_{j,l}) \geq h_0^*(x_{j,l}) + \gamma_{i,l} h_0^*(x_{i,l}), \quad \forall i \in \{1, \dots, k\} \forall j \in \{k+1, \dots, N\}. \quad (26)$$

In the following computation, we take a random patient $i^* \in \{1, \dots, k\}$ and a random patient $j^* \in \{k+1, \dots, N\}$ and exchange the treatment—that is, we provide preventive treatment to patient j^* but not to patient i^* . From this exchange, we arrive at the following:

$$n_{\text{onsets}} = \sum_{\substack{i=1 \\ i \neq i^*}}^k \mathbb{E}[y | x_{i,l}, t_{i,l} = 1] + \sum_{\substack{j=k+1 \\ j \neq j^*}}^N \mathbb{E}[y | x_{j,l}, t_{j,l} = 0] + \mathbb{E}[y | x_{j^*,l^*}, t_{j^*,l^*} = 1] + \mathbb{E}[y | x_{i^*,l^*}, t_{i^*,l^*} = 0] \quad (27)$$

$$= \sum_{\substack{i=1 \\ i \neq i^*}}^k (1 - \gamma_{i,l}) h_0^*(x_{i,l}) + \sum_{\substack{j=k+1 \\ j \neq j^*}}^N h_0^*(x_{j,l}) + (1 - \gamma_{j^*,l^*}) h_0^*(x_{j^*,l^*}) + h_0^*(x_{i^*,l^*}) \quad (28)$$

$$\geq \sum_{\substack{i=1 \\ i \neq i^*}}^k (1 - \gamma_{i,l}) h_0^*(x_{i,l}) + \sum_{\substack{j=k+1 \\ j \neq j^*}}^N h_0^*(x_{j,l}) + (1 - \gamma_{i^*,l^*}) h_0^*(x_{i^*,l^*}) + h_0^*(x_{j^*,l^*}) \quad (29)$$

$$= \sum_{i=1}^k (1 - \gamma_{i,l}) h_0^*(x_{i,l}) + \sum_{j=k+1}^N h_0^*(x_{j,l}) \quad (30)$$

$$= \sum_{i=1}^k \mathbb{E}[y | x_{i,l}, t_{i,l} = 1] + \sum_{j=k+1}^N \mathbb{E}[y | x_{j,l}, t_{j,l} = 0] \quad (31)$$

Thus, exchanging the treatments of patients i^* and j^* leads to the same or a larger number of expected onsets.

Consequently, our original allocation minimizes the number of expected onsets. \square

Appendix C: Robustness Checks

C.1. Choice of Machine Learning Model

For the analysis in the main paper, the machine learning in stage 2 of the decision model was set to gradient boosted decision trees (Friedman 2001). Gradient boosted decision trees belong to the category of tree ensemble methods, which have been applied successfully in other operational applications (e. g., Glaeser et al. 2019, Senoner et al. 2022). We now provide empirical evidence supporting this choice. For this purpose, we compare the performance of our decision model under different machine learning models in stage 2 while the rest of the decision model remains the same. Our comparison includes linear (lasso and ridge regression) and nonlinear (random forest and deep neural network) models. All estimation details are provided in Supplement G.

Results are shown in Table 7. While the lasso performs well among the linear models, nonlinear machine learning models consistently outperform linear models. This may be because nonlinear models are better able to exploit complex patterns in the data, particularly with regard to between-patient heterogeneity. However, nonlinear models can also be prone to overfitting, which may explain why the deep neural network underperformed compared to the gradient boosted decision trees. Overall, the decision model using gradient boosted decision trees achieves the best performance.

Table 7 Performance comparison among different machine learning models.

	$k=1,000$		$k=5,000$		$k=10,000$	
	Prevented diseases	Cost savings	Prevented diseases	Cost savings	Prevented diseases	Cost savings
Naïve baseline	53.438 (0.330)	6.319 (0.235)	250.422 (1.698)	8.889 (0.258)	489.331 (2.162)	11.988 (0.249)
Lasso	152.124 (0.840)	7.489 (0.218)	520.490 (2.401)	12.726 (0.234)	811.841 (3.684)	16.800 (0.261)
Ridge regression	150.495 (0.831)	7.409 (0.216)	514.812 (2.374)	12.586 (0.231)	803.021 (3.644)	16.617 (0.258)
Random forest	160.270 (0.885)	7.892 (0.230)	548.875 (2.531)	13.425 (0.246)	855.942 (3.880)	17.717 (0.276)
Deep neural network	155.636 (0.860)	7.670 (0.224)	535.133 (2.467)	13.109 (0.240)	833.792 (3.766)	17.276 (0.268)
Gradient boosted decision trees	162.913 (0.902)	8.050 (0.239)	567.710 (2.612)	13.979 (0.253)	882.004 (3.934)	18.337 (0.284)

Stated: mean performance (standard deviation in parentheses)

Note. Performance metrics for allocation of preventive care when using the respective machine learning model in stage 2 of our decision model. Budget constraints allow for the treatment of k patients per year. Cost savings are reported in USD millions.

C.2. Comparison with Alternative Risk Scores

The Framingham diabetes risk score represents the quasi-standard in clinical practice for assessing diabetes risk (Long and Fox 2016). Nevertheless, Table 8 shows the comparison between our decision model and other diabetes

risk scores from clinical practice—namely, Lindstroem and Tuomilehto (2003), Wilson et al. (2007), Kahn et al. (2009), and Rosella et al. (2011). Consistent with our previous findings, our proposed decision model outperforms all alternative risk scores.

Table 8 Performance comparison between our decision model and other risk scores.

	$k=1,000$		$k=5,000$		$k=10,000$	
	Prevented diseases	Cost savings	Prevented diseases	Cost savings	Prevented diseases	Cost savings
Naïve baseline	53.438 (0.330)	6.319 (0.235)	250.422 (1.698)	8.889 (0.258)	489.331 (2.162)	11.988 (0.249)
Lindstroem and Tuomilehto (2003)	128.810 (0.710)	6.323 (0.181)	434.431 (2.008)	10.562 (0.196)	679.740 (3.126)	14.015 (0.219)
Wilson et al. (2007)	121.424 (0.671)	5.979 (0.174)	415.914 (1.918)	10.173 (0.187)	648.571 (2.940)	13.425 (0.209)
Kahn et al. (2009)	119.542 (0.660)	5.879 (0.170)	406.948 (1.878)	9.930 (0.183)	635.441 (2.897)	13.133 (0.205)
Rosella et al. (2011)	120.917 (0.667)	5.940 (0.171)	409.335 (1.891)	9.966 (0.184)	639.950 (2.932)	13.207 (0.206)
Our decision model	162.913 (0.902)	8.050 (0.239)	567.710 (2.612)	13.979 (0.253)	882.004 (3.934)	18.337 (0.284)

Stated: mean performance (standard deviation in parentheses)

Note. Performance metrics for allocation of preventive care with the respective decision model. Budget constraints allow for the treatment of k patients per year. Cost savings over no preventive care allocation are reported in USD millions.

C.3. Robustness to Estimation Errors in the Treatment Effect

We now analyze the robustness of our decision model to potential errors in the estimation of the treatment effect. For this, we run the same analysis as before, but introduce additional Gaussian noise to the estimated treatment effect $\gamma_{i,l}$, i. e.,

$$\tilde{\gamma}_{i,l} = \gamma_{i,l} + \epsilon \quad \text{with} \quad \epsilon \sim \mathcal{N}(0, \sigma^2), \quad (32)$$

where σ denotes the level of noise. The latter thus represents the estimation error. We further know that the effectiveness of preventive care such as *metformin* is bounded and, therefore, ensure the noisy treatment effect $\tilde{\gamma}_{i,l}$ to lie in the interval $[0, 1]$ through clipping.

Table 9 lists the results for this analysis. Here, we compare different noise levels $\sigma \in \{0, 0.1, 0.5\}$. Of note, a noise level of $\sigma = 0.5$ introduces a large estimation error, especially when considering that the effectiveness of preventive treatments should lie in the interval $[0, 1]$ for real-world clinical settings. As expected, we find a slightly more dominant role of estimation errors for small k , which, in light of the small set of patients allocated to preventive care, results in slightly reduced cost savings. This can be expected due to the small sample size. However, and more importantly, we find that the impact of estimation errors is overall fairly small: we still achieve a large number of prevented disease onsets as well as large cost savings over current practice. For example, even

for noise of $\sigma = 0.1$, our decision model still outperforms current practice for $k = 10,000$ and is only slightly outperformed by current practice for noise of $\sigma = 0.5$.

Table 9 Robustness of our decision model to estimation errors in the treatment effect.

	$k=1,000$		$k=5,000$		$k=10,000$	
	Prevented diseases	Cost savings	Prevented diseases	Cost savings	Prevented diseases	Cost savings
Clinical baseline	134.798	7.340	446.948	11.996	701.942	15.350
$\sigma = 0$	(1.025)	(0.245)	(3.373)	(0.258)	(3.837)	(0.265)
$\sigma = 0.1$	162.913	8.050	567.710	13.979	882.004	18.337
	(0.902)	(0.239)	(2.612)	(0.253)	(3.934)	(0.284)
$\sigma = 0.1$	149.481	7.331	501.655	12.172	785.779	16.180
	(0.823)	(0.209)	(2.320)	(0.227)	(3.630)	(0.253)
$\sigma = 0.5$	131.199	6.462	449.977	11.012	701.492	14.526
	(0.725)	(0.189)	(2.075)	(0.202)	(3.176)	(0.226)

Stated: mean performance (standard deviation in parentheses)

Note. Budget constraints allow for the treatment of k patients per year. Cost savings over no preventive care allocation are reported in USD millions.

C.4. Robustness to Number of Training Samples

In machine learning, the size of the training set is known to affect the prediction performance, which, in turn, affects the operational performance. We thus study the convergence of our decision model against the optimal decision for an increasing number of observations. We use a static setting with observations $x \in \mathbb{R}^3$, for which the three components represent different patient variables (age, height, body mass index). We simulate these as follows:

$$x \sim \mathcal{N}(\mu, \Sigma), \quad (33)$$

$$\mu = (50, 170, 27)^T, \quad (34)$$

$$\Sigma = \begin{bmatrix} 20 & 0 & 5 \\ 0 & 50 & 5 \\ 5 & 5 & 5 \end{bmatrix} \quad (35)$$

The true but unobservable risk of an onset is given by

$$y = \sigma \left(\frac{0.5x_1 + 0.1x_2 + 0.2x_3}{100} \right), \quad (36)$$

with σ denoting the sigmoid function. We assume that we observe the actual outcomes only for N_{train} patients.

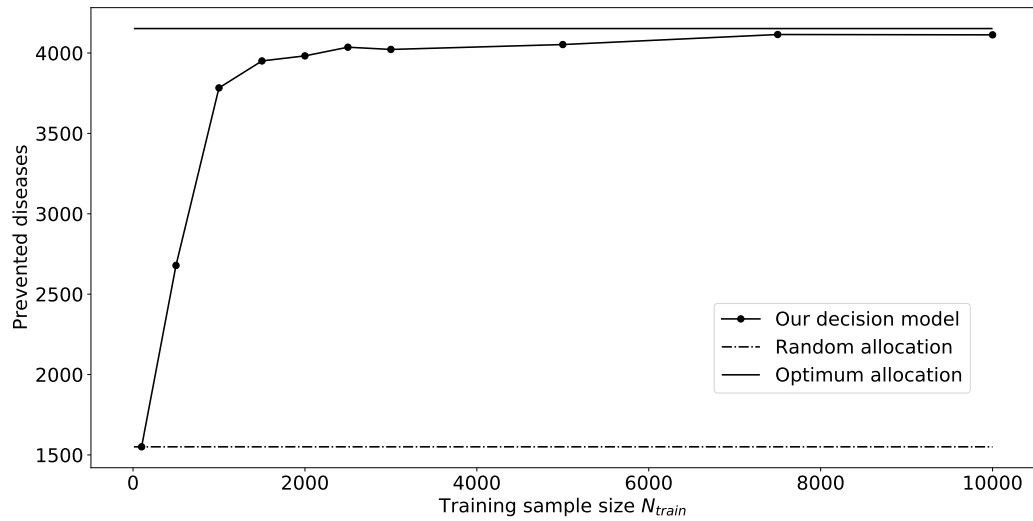
That is, for each patient i , we assume that we observe the realization of Y given by

$$y_i = \begin{cases} 1, & \text{if } y + \epsilon \geq 0.7, \\ 0, & \text{otherwise,} \end{cases} \quad \text{with } \epsilon \sim \mathcal{N}(0, 1). \quad (37)$$

In the simulation, we further set the true treatment effect to $\gamma = 0.31$.

Our population comprises $N = 100,000$ patients. The budget constraint allows us to enroll $k = 10,000$ patients into preventive care. Figure 4 shows the number of prevented diseases for varying training sample sizes N_{train} . It confirms that our decision model converges to the optimal allocation. Moreover, we already yield a close-to-optimum performance for a relatively small number of training samples (i. e., fewer than 3,000).

Figure 4 Performance of our decision model for varying number of training samples.



Appendix D: Causal Sensitivity Analysis

In the following, we analyze how potential unobserved confounders would affect the estimated treatment effects. Unobserved confounders refer to unmeasured variables that affect both treatment assignment and the health outcome (Cinelli and Hazlett 2020). A prominent example is the socioeconomic status of patients, as wealthier patients are often more likely to enroll in preventive care and, furthermore, can afford a healthier lifestyle that is responsible for a lower diabetes risk, which thus may underestimate the treatment effect of *metformin*. Since these confounders are unobserved, we cannot control for them when estimating the effect of a treatment. As a remedy, we follow the so-called causal sensitivity analysis from Cinelli and Hazlett (2020) and analyze the effect that an unobserved confounder would have on the treatment effect estimation, if it existed. More precisely, we select an observed variable which affects both treatment assignment and the health outcome and now assume that an unobserved confounder exists with the same strength as the observed one. Then, we check that the treatment effect cannot be explained away and thus remains robust.

The age of a patient is generally associated with the likelihood of being prescribed with metformin (e.g., Rosella et al. 2011) and has also been shown to be strongly associated with the risk of diabetes onset (e.g., Knowler et al. 2002). Thus, we use the strength of the variable “age” in the following as our reference to make comparisons of whether the treatment effect can be explained away.

We follow Cinelli and Hazlett (2020) and use a linear regression. Informed by clinical research (e.g., the Framingham risk score), we use a linear regression with the following patient variables, i.e.,

$$y = treatment + age + sex + weight + body_mass_index + systolic_bp + diastolic_bp + height. \quad (38)$$

We then perform causal sensitivity analysis for all patients and first test for the sensitivity of the main effect. Afterward, to account for heterogeneity in the treatment effect, we repeat the causal sensitivity analysis for different patient subgroups. For this, we segment the patients by their estimated treatment effects into four groups – A, B, C, D – using a decision tree regressor, and then perform the causal sensitivity analysis for each subgroup.

Table 10 shows the results of our causal sensitivity analysis. The results confirm that the estimated treatment effect is robust to unobserved confounders. The results are also consistent across all subgroups. Note that the patient’s age is considered to be one of the strongest predictors for diabetes (American Diabetes Association 2022), and, thus, it is very unlikely to have an unobserved confounder of a similar strength. In other words,

the treatment effect cannot be explained away by unobserved confounding. In sum, we do not see evidence that unobserved confounders may undermine our estimates of treatment effectiveness, which is in line with many existing works using randomized control trials that have already confirmed the effectiveness of *metformin* (e.g., Knowler et al. 2002).

Table 10 Results of robustness check for unobserved confounders.

Strength	All patients	Subgroups			
		A	B	C	D
$0.2 \times age$	0.066 [0.052, 0.088]	0.066 [0.045, 0.096]	0.142 [0.103, 0.188]	0.048 [0.026, 0.074]	0.058 [0.027, 0.097]
$0.5 \times age$	0.053 [0.039, 0.074]	0.052 [0.029, 0.078]	0.127 [0.084, 0.177]	0.037 [0.015, 0.060]	0.046 [0.019, 0.075]
$0.8 \times age$	0.043 [0.031, 0.058]	0.040 [0.013, 0.072]	0.117 [0.075, 0.132]	0.029 [0.009, 0.052]	0.039 [0.014, 0.071]

Stated: mean estimation (95% confidence intervals in parentheses)

Note. Estimated treatment effects when assuming the existence of an unobserved confounder with the strength $0.2 \times age$, $0.5 \times age$, and $1.0 \times age$.

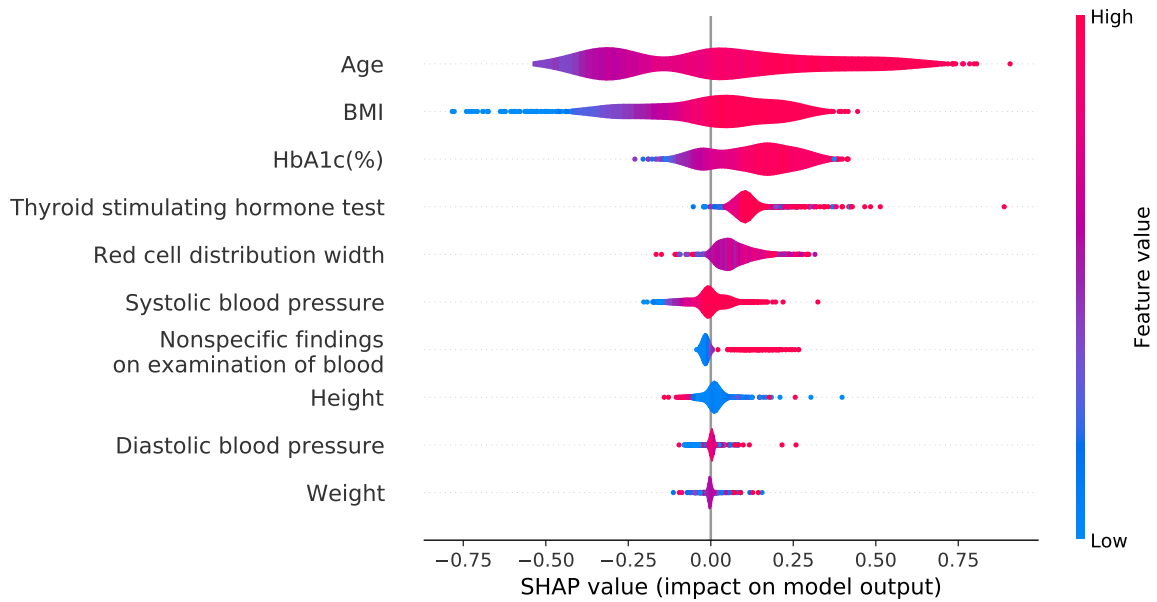
Appendix E: Machine Learning Explainability

To determine the most influential predictors, we calculated SHAP values. SHAP values are a unified method for measuring how much a predictor contributes to the overall model output and, therefore, rank the importance of the features. Moreover, they indicate whether higher (lower) values of a predictor are associated with an increased risk of diabetes onset.

Figure 5 shows the ten most important predictors for the gradient boosted decision tree when predicting diabetes onset. The most important predictors are age, body mass index (BMI), and the current glycated hemoglobin level (HbA1c). Age and BMI are well known risk factors for diabetes, and both were accordingly among the most important predictors for our gradient boosted decision tree (Abbasi et al. 2016). Also among the ten most important predictors are the red cell distribution width and thyroid stimulating hormone tests, which have both been found in the medical literature to be predictors for diabetes (e.g., Chaker et al. 2016). The height and weight of the patient are also among the ten most important predictors, yet with less importance. This is likely due to the fact that the predictive power of these variables as individual predictors is limited, and only their combination (as in the BMI) leads to an important predictor.

We discussed our results with medical experts from diabetes care, who confirmed to us that these predictors are well-established risk factor for diabetes.

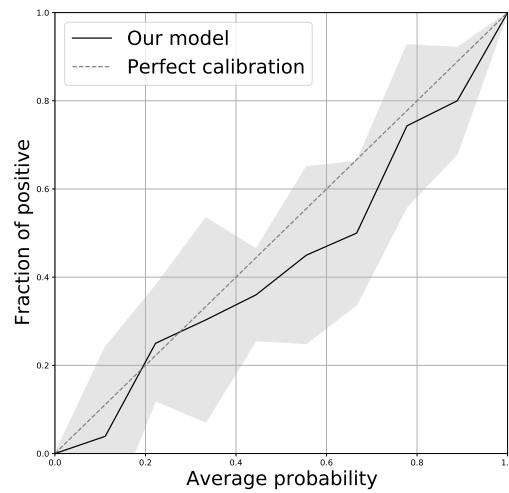
Figure 5 SHAP values for the ten most important predictors for the prediction model. The dots represent observation for which the x-axis denotes the SHAP value. Observations with a positive (negative) SHAP value denote a higher (lower) risk for a diabetes onset. The color of each dot indicates the corresponding value of the predictor for a given observation.



Appendix F: Calibration of Machine Learning Model

We use two techniques to improve the calibration of our machine learning model. Here, the aim is that the distribution of the predicted probability is similar to the distribution of the observed probability in the training data. First, we use the synthetic minority over-sampling technique (SMOTE) (Chawla et al. 2002) to oversample observations from patients who develop diabetes. Second, we use Platt scaling to calibrate the predicted class probabilities (Niculescu-Mizil and Caruana 2005). Figure 6 shows the calibration plot for our gradient boosted decision trees after resampling. The calibration plot shows a good calibration.

Figure 6 Calibration of our machine learning model. The gray area shows the standard deviation due to resampling.



Appendix G: Hyperparameter Tuning

We determined the hyperparameters using Bayesian optimization with 100 iterations of optimization. Table 11 lists the tuning ranges for each parameter. In the Bayesian optimization, the initial parameters are sampled from a uniform distribution. We implemented and evaluated all decision models in Python. The lasso, ridge regression, and random forest machine learning models are implemented using scikit-learn. The deep neural network is implemented using Pytorch. We train the deep neural network using the Adam optimizer. We use ReLU activation in the hidden layers and sigmoid in the output layer. The objective of the optimizer is to minimize the binary cross-entropy. The gradient boosted decision trees are implemented using LightGBM. For the causal forest, we keep all parameters at their default values. That is, we set the number of estimators to 10, the number of features to consider when looking for the best split to 10, the maximum depth of the trees to 5, the minimum number of samples required to be split at a leaf node to 100, and the minimum number of samples required of the experiment group to be split at a leaf node to 10.

Table 11 Hyperparameter tuning.

Model	Tuning parameters	Tuning range
Lasso	Regularization strength α	(0.001, 10000)
Ridge regression	Regularization strength	(0.001, 10000)
Random forest	Number of trees	(20, 200)
	Min samples for split	(2, 150)
	Number of features to consider for split	{sqrt, log2, all}
	Criterion for split	{gini, entropy}
	Class weight	{none, balanced}
Deep neural network	Number of neurons	(5, 128)
	Number of hidden layers	(1, 5)
	Regularization strength	(0.0, 5.0)
	Learning rate	(0.0001, 0.1)
	Batch size	(10, 1024)
Gradient boosted decision trees	Number of trees	(20, 200)
	Number of leaves	(20, 150)
	Learning rate	(0.01, 0.5)
	Number of samples for constructing bins	(20000, 300000)
	Minimum number of data needed in a child	(20, 500)
	L1 regularization term on weights	(0, 1)
	L2 regularization term on weights	(0, 1)
	Class weight	{none, balanced}

Appendix H: Applicability for Settings with Known Treatment Effects

Here, we aim to demonstrate the applicability of our model in settings where the treatment effect of preventive care is known. In such a case, there is no need for healthcare organizations to perform counterfactual inference for treatment effect estimation, as the treatment effect can be directly entered in the decision model.

To show the applicability of our decision model, we build upon a different setting. Recall that we examined *metformin* in the main paper. The reason was that it represents the quasi-standard for preventive care aimed at individuals at risk of developing diabetes in current medical practice (Long and Fox 2016). Several studies also suggest a risk reduction if patients are enrolled in personalized coaching toward a healthier lifestyle; however, such personalized coaching is currently not covered by health insurers in many countries (e. g., Germany, Switzerland, and Medicare in the United States). Nevertheless, it may be interesting to study the cost-effectiveness of lifestyle coaching for diabetes prevention.

In the case of personalized coaching, the goal of preventive treatment is to achieve lifestyle changes of the patients. Lifestyle changes have not been prescribed to the customers of our partnering health insurer. Thus, the treatment effect cannot be estimated through stage 1 of our decision model. Instead, we rely on prior literature that has estimated the treatment effect of personalized coaching through randomized controlled trials (Knowler et al. 2002).

We assume a treatment effect due to personalized coaching of $\gamma_{i,l^*} = 0.58$ (Knowler et al. 2002) that is exponentially decreasing over time, i. e.,

$$\gamma_{i,l} = \gamma_{i,l^*+1} \exp(l^* - l), \quad \text{for } l \geq l^*. \quad (39)$$

We set the annual cost of personalized coaching to $C_{\text{prevent}} = 1,600$ (Azelton et al. 2021).

Table 12 shows the empirical results. Here, we again compare the data-driven allocations from our decision model against that from current practice. Our findings are line with the results from the main paper: our decision model outperforms current practice across all metrics.

To sum up, there are different benefits of whether the treatment effect is (a) known or (b) estimated through stage 1 of our decision model. A benefit of (a) is that healthcare organizations can leverage known treatment effects from randomized controlled trials, which represent the gold standard for measuring treatment effects. In particular, there is no bias due to unobserved confounders. A benefit of (b) is that randomized controlled trials are not always available or otherwise costly, because of which our estimation may be preferred in practice. Moreover, using counterfactual inference in our decision model offers a mathematical approach to directly

learn heterogeneous treatment effects and thus account for the differential effectiveness of preventive care across patients.

Table 12 Comparison of current practice vs. our data-driven decision model for the case of personalized coaching.

	$k=1,000$		$k=5,000$		$k=10,000$	
	Prevented diseases	Cost savings	Prevented diseases	Cost savings	Prevented diseases	Cost savings
Naïve baseline	728.378 (5.606)	16.457 (0.141)	1423.607 (10.761)	26.651 (0.249)	2736.326 (16.553)	45.807 (0.323)
Current practice	1318.303 (9.893)	25.786 (0.223)	2083.693 (12.719)	37.088 (0.287)	3184.516 (14.837)	51.263 (0.310)
Our decision model	1555.792 (10.960)	29.544 (0.223)	2143.867 (14.553)	38.036 (0.295)	3349.127 (15.103)	53.767 (0.276)

Stated: mean performance (standard deviation in parentheses)

Note. Performance metrics for allocating preventive care given a varying budget for enrolling k patients per year into preventive treatments. Cost savings over no preventive care allocation are reported in USD millions.