

Treatment effect estimation by comparing observed and predicted outcomes: theory and practical illustration in radiotherapy

Lotta M. Meijerink
Bas B. L. Penning de Vries
Remi A. Nout

Artuur M. Leeuwenberg
Johannes A. Langendijk
Karel G.M. Moons

Jungyeon Choi
Judith G.M. van Loon
Ewoud Schuit

November 27, 2025

Abstract

Prediction models developed before the introduction of a new treatment may be used to estimate treatment effects of newly introduced treatments. One approach, known as model-based clinical evaluation in radiotherapy, does this by comparing *observed* outcomes under a new treatment with *predicted* outcomes had these patients received the standard treatment. This article clarifies the relevant conditions needed for valid average treatment effect estimation using this approach, using the potential outcomes framework and a practical case study.

1 Introduction

New treatments need extensive evaluation to confirm their (cost-)effectiveness and safety. While randomized controlled trials (RCTs) are considered the gold standard for estimating treatment effects [1], they are not always feasible or ethical, or may not yet be available at the time decisions about the new treatment need to be made [2, 3]. Consequently, treatments are sometimes introduced before trial evidence is available [4]. Also, there may be interest in rare outcomes or patient populations not covered by previous trials. In such cases, non-randomized methods can be used to generate evidence on treatment effectiveness, either while awaiting trial results, as a substitute for a trial, or to complement existing evidence. Commonly used methods include using regression adjustment or propensity score techniques.

An alternative approach focuses on the set of patients who received the newly introduced treatment and applies a model to predict their counterfactual outcomes: *what would their outcome have been had they received the standard treatment instead?* The average treatment effect among the treated (ATT) is estimated as the average difference between observed outcomes under the new treatment and predicted outcomes under the standard treatment. In this approach, the model used for counterfactual predictions is based on historical or external patient data from a setting where the new treatment has not yet been implemented, i.e., where everyone received the standard treatment. The approach can be recognized as a specific case of the parametric g-formula or g-computation [5], or an application of model-based standardization [6, 7].

In radiotherapy, it was proposed as model-based clinical evaluation, or a model-based approach [4, 8]. There it is used to estimate the effect of new technologies, such as proton therapy or swallowing sparing intensity modulated radiotherapy, in reducing radiation-induced complications [9–11].

While the approach (comparing observed outcomes under the new treatment with predicted outcomes under the standard treatment) seems intuitive and practical, it relies on several conditions that are often not explicitly discussed. The aim of this paper is to clarify these conditions. We do this by formalizing the approach in the potential outcomes framework and discussing a set of sufficient conditions: what they mean, why we need them, and when they would be violated. To illustrate the approach and its conditions more concretely, we discuss them using a realistic case study, which is introduced below.

2 Case study

2.1 Background

Radiation technology is constantly developing, often aiming to reduce the radiation dose to healthy organs surrounding the tumor, and thereby preventing radiation-induced complications, while maintaining the required dose to the tumor itself. An example is proton therapy, which allows for more precise radiation delivery - specifically, less dose to surrounding healthy organs - compared to ‘conventional’ photon therapy (such as VMAT: Volumetric Modulated Arc Therapy). Clinicians widely agree that in general, reducing radiation exposure to healthy tissues is beneficial, following the ALARA principle that radiation exposure should be “as low as reasonably achievable”.

As a result, innovative radiotherapy technologies are often adopted relatively quickly, sometimes even without formal evaluation through RCTs [4]. Furthermore, there are additional barriers to conducting RCTs on new radiation technologies. For example, unlike pharmaceutical companies, medtech firms have limited financial incentives to fund trials, as successful technologies generally benefit all providers rather than just the one investing in the often costly study. Nonetheless, especially when investment costs are considerable, gathering evidence of the benefits of the new technologies remains important [10].

2.2 Research aim

This case study is based on (but a simplified version of) a real-world setting in the Netherlands and focuses on patients with head and neck cancer undergoing radiotherapy. All example calculations are based on synthetic data. Interest was in the estimation of the benefit of proton therapy in terms of reducing dysphagia (difficulty swallowing) at 6 months after radiotherapy, compared to photon therapy (VMAT), or more specifically, in the average benefit in the patients currently eligible for proton therapy. Here, eligibility was determined by ‘model-based selection’ [12], which is further explained below.

2.3 Population

Patients were treated in a single center, covering two time periods:

- **2007-2017 (Pre-introduction sample):** All 750 patients treated during this period received photon-based radiotherapy (VMAT).
- **2018-2019 (Post-introduction sample):** During this time, proton therapy became available as a treatment option, but due to its higher costs and limited capacity, was only reimbursed for patients who were expected to benefit significantly from it. Not all patients are expected to benefit equally from proton therapy. For example, in some cases, the tumor is located such that conventional photon therapy can target it effectively, without causing much damage to nearby healthy tissues. In such situations, the added value of the more expensive proton therapy is minimal. However, for other patients, proton therapy can substantially reduce the radiation dose delivered to surrounding healthy organs and thereby reduce radiation-induced toxicity.

A model-based selection approach was introduced to determine which patients would be eligible to receive proton therapy.[12] In this approach, for each patient, two radiation plans were created: one using photon therapy and one using proton therapy. From these plans, mean planned doses to tissues surrounding the tumor (the superior, middle, and inferior pharyngeal constrictor muscles and the oral cavity) were extracted. These dose metrics, along with two patient-specific characteristics, baseline dysphagia status and tumor location, served as input for an existing logistic regression model to predict the risk of dysphagia at 6 months after radiotherapy. The difference in predicted dysphagia risk between the photon and proton plans was used as an estimate of the expected benefit of proton therapy. Patients were selected for proton therapy if this expected benefit exceeded 10%. To summarize, selection for proton therapy was based directly on the planned dose parameters from both treatment plans and the two baseline characteristics. Using this procedure, 93 patients were selected for and treated with proton therapy, while 207 received photon therapy during this period.

For all patients, data were collected on baseline characteristics, the four photon plan dose parameters, and the grade of dysphagia after 6 months (the outcome). For patients treated after the introduction of proton therapy, additionally, the same four dose parameters were available from the proton plan. We refer to Appendix A.1 for a synthetic example of what the data of this case study could look like, and Appendix B for the description of how these synthetic data were generated.

3 Methodology

This section describes the methodology in general, and in the context of the case study. We start with describing what we are trying to estimate (the estimand), followed by how this is estimated, the sufficient conditions, and finally, ways to gain supportive evidence of the validity of the approach.

3.1 Causal estimand

The approach is used to compare two treatments:

- $T = 0$: the standard (control) treatment
- $T = 1$: the (new) target treatment under evaluation

Let $Y \in \{0, 1\}$ be a binary outcome, where $Y = 1$ indicates the presence or occurrence of a specific health state at the predefined relevant outcome timing. Our causal question is:

To what extent does the target treatment $T = 1$ reduce the risk of the outcome $Y = 1$, compared to standard treatment $T = 0$, in the population currently treated with $T = 1$?

This corresponds to the average treatment effect among the treated (ATT). To formalize this question, we define it using the potential outcomes framework [13]. We let $Y(t)$ denote the potential outcome of an individual under treatment $T = t$, i.e., the outcome that would be observed if the individual receives treatment t . Then, the ATT is defined as:

$$ATT = \mathbb{E}_{post}[Y(1) - Y(0) \mid T = 1]$$

i.e., the average difference between the two potential outcomes, in the patients treated with the target therapy. Here, we use the subscript *post* to denote the ‘post-introduction’ setting, i.e. the setting where the target treatment is already introduced. This setting will be different from the setting that provides us data to develop the model (e.g. because they come from a different hospital or moment in time), which we will refer to as ‘pre-introduction’.

Alternative effect measures (e.g., odds ratio or risk ratio) can be defined analogously.

Case study In the case study, photon therapy is the standard treatment $T = 0$, and proton therapy is the target treatment $T = 1$, and the outcome Y we consider is the presence of dysphagia (grade 2 or higher) at 6 months after radiotherapy. For an individual, the potential outcome $Y(0)$ represents whether they would experience dysphagia if they receive photon therapy ($T = 0$). The population of interest is the population that is currently treated with protons, and we want to know to what extent proton therapy reduces the risk of dysphagia, compared to photon therapy, in that population. Or in other words, the expected difference in proportion of patients with dysphagia if everyone in that population were treated with photons versus if everyone were treated with protons. We care about the treatment effect in the population currently treated with protons because these are the patients in whom we expected a benefit of proton therapy.

3.2 Estimation

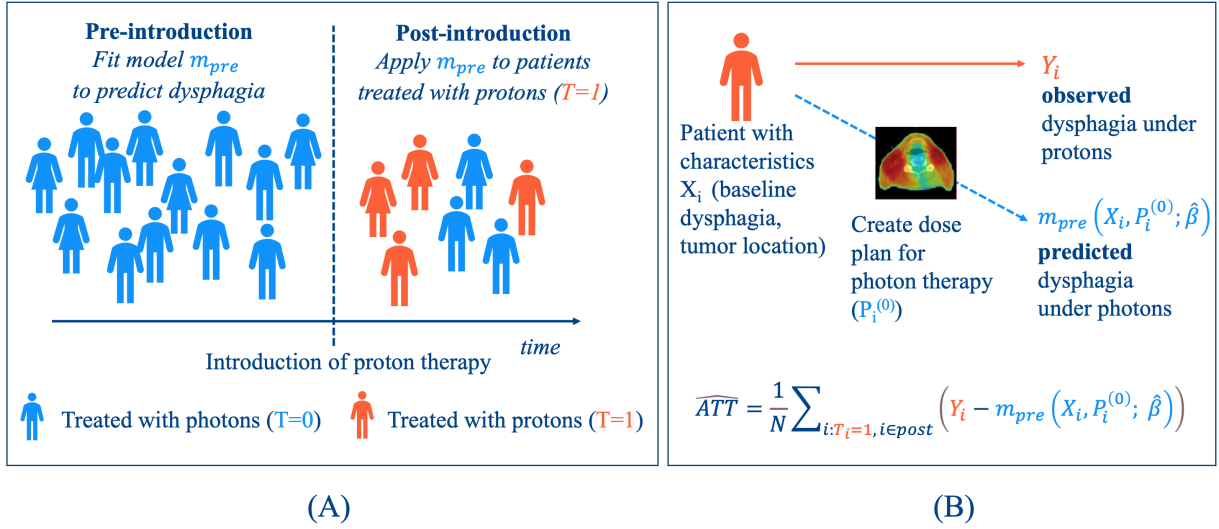


Figure 1: Illustration of the approach as applied in the case study. (A) The patients that were treated before introduction of proton therapy were used to develop the model \mathbf{m}_{pre} . These patients were all treated with photons (represented in blue). After introduction of proton therapy, only a subset was treated with protons (represented in orange). (B) On this proton-treated subset of individuals, the model is applied to make counterfactual predictions using the previously developed model \mathbf{m}_{pre} . To make predictions, the model uses four variables from their individual photon dose plan $\mathbf{P}^{(0)}$ and their other patient characteristics \mathbf{X} (baseline dysphagia and tumor location). The average difference between predicted and observed outcomes is taken as an estimate of the average treatment effect among the treated (ATT).

To estimate the ATT, we make use of a model family $\mathbf{m}_{pre}(\mathbf{X}, \mathbf{P}^{(0)}, \beta)$, parametrised by β (e.g. a logistic regression with a set of coefficients). We fit the model in a historical (or external) cohort where all patients received the standard treatment $T = 0$, giving us fitted parameters $\hat{\beta}$. This fitted model estimates the risk of the outcome $Y = 1$ under treatment $T = 0$ given:

- **Patient characteristics \mathbf{X}** , e.g., age, comorbidities
- **Treatment plan variables $\mathbf{P}^{(0)}$ of the standard treatment**, such as planned radiation doses to specific organs under photon therapy in the case study. In case the standard treatment of interest can only have one single operationalization (e.g., a pill with fixed dose), there will be no need to include treatment plan variables.

Let $i = 1, \dots, N$ index the sample of patients from the post-introduction population, that are currently treated with $T = 1$. Then:

$$\widehat{ATT} = \frac{1}{N} \sum_{i=1}^N \left(Y_i - m_{pre}(\mathbf{X}_i, \mathbf{P}_i^{(0)}; \hat{\beta}) \right)$$

In words, this estimate is the average difference between observed outcomes and predicted counterfactuals under $T = 0$, in the sample of the post-introduction population that was treated with $T = 1$.

Case study The estimation procedure of the case study is visualized in Figure 1. The model family \mathbf{m}_{pre} is a logistic regression model, and its parameters β are fitted on the 750 individuals treated with photon therapy between 2007 and 2017. As patient characteristics \mathbf{X} , the model includes the grade of dysphagia at baseline (i.e., before starting radiotherapy) and tumor location. As treatment plan characteristics $\mathbf{P}^{(0)}$, four variables from the photon radiation plan are used: the mean planned dose to the superior pharyngeal constrictor muscle (PCM), middle PCM, and inferior PCM. The ATT is calculated using the 93 individuals treated with protons in the target sample (2018?2019). For these individuals, although they were treated with protons, corresponding photon plans are also available as these were used selecting patients for proton therapy (see section 2.3). This allows us to compare their observed dysphagia outcome Y_i with the counterfactual risk of dysphagia predicted by the fitted logistic regression model, using their photon plan and patient characteristics: $\mathbf{m}_{pre}(\mathbf{X}_i, \mathbf{P}_i^{(0)}; \hat{\beta})$. We refer to Appendix A.4 for an example calculation in R, using bootstrapping to obtain confidence intervals around the estimated ATT.

3.3 Conditions for validity

As illustrated above, it is relatively easy to estimate a treatment effect using the approach above, but under what conditions does it provide a valid answer to the causal question? We describe a set of conditions that ensure validity of the estimation procedure. Weaker conditions could suffice to obtain valid estimates; however, we present the current set of conditions because they offer an intuitive and useful framework for understanding the problem. For more details, including a derivation from causal estimand to the estimator, that shows why these conditions are sufficient, please refer to Appendix C.

3.3.1 Condition 1: Transportability

This condition requires that the risk of outcome Y under treatment $T = 0$ *given covariates in the model* are the same in the pre- and post-introduction populations:

$$P_{pre}(Y(0) = 1 \mid \mathbf{X}, \mathbf{P}^{(0)}) = P_{post}(Y(0) = 1 \mid \mathbf{X}, \mathbf{P}^{(0)})$$

This holds if any variable influencing outcome $Y(0)$ has the same distribution across populations (conditional on covariates) or is correctly incorporated in the model. In general, the relationship between covariates and the outcome also must stay consistent between populations. The condition only applies to $Y(0)$, the outcomes under treatment $T = 0$, as only these are predicted using the pre-introduction model.

Why? This condition allows us to apply a model developed in a sample from a historical or external population to make valid counterfactual predictions in the post-introduction population.

Examples of violations

- If smoking status is not included in the model but does increase the outcome risk *and* there is a higher smoking prevalence in the pre-introduction population, compared to post-introduction. This would cause the model to overestimate risk when applied to the post-introduction population.
- If there are concurrent changes in treatment protocols (e.g., different opioid prescribing protocols between hospitals or over time) that influence the outcome but aren't included as predictors in the model.
- If there are differences in dose optimization strategies between populations, causing differences in the radiation dose to the larynx, and this dose does affect the outcome, but is not incorporated in the model.

Case study As described in section 2.3, the case study considers a situation where pre- and post-introduction samples came from the same center, so potential differences would primarily have arisen from temporal changes (in contrast to a situation where the pre-introduction sample used to develop the model was collected in a different hospital, for example). There were no notable changes regarding concurrent treatment strategies and changes in patient characteristics are small, as only a relatively short

time span is considered. However, over time, dose optimization strategies improved, leading to reduced doses in nearly all organs, including those not explicitly included in the model. While the organs most relevant to the outcome of interest (dysphagia) are represented in the model, these temporal changes may still cause an *overestimation* of predicted outcomes under photon therapy. This would result in an *overestimation* of the treatment effect, meaning we might detect a larger benefit of protons than truly exists, due to the improvement of the dose planning under the standard treatment, i.e., photon therapy.

Another important consideration that may influence the relation between covariates and the outcome, is whether in the population treated with protons, the *photon plans* accurately represent what their photon plan would have been had they received photon therapy, i.e. that the quality of the photon plan is not impacted by the fact that the patient is selected for proton therapy. In the case study, it is likely that these plans are indeed realistic because these plans were created before treatment allocation decisions (proton vs. photon therapy) were made, making it unethical and unrealistic that poor photon plans were created.

3.3.2 Condition 2: Ignorability of treatment assignment

This condition requires that within the post-introduction population, the potential outcome that would be observed under standard treatment is independent of the treatment that the patient was actually treated with, given covariates:

$$Y(0) \perp T \mid \mathbf{X}, \mathbf{P}^{(0)}$$

Why? This condition allows us to apply a model developed in a sample from a historical or external population, where *everyone* received the standard treatment, to make valid counterfactual predictions in a non-random *subset* of the post-introduction population: the individuals that were selected to be treated with the target treatment. If selection was random, this condition is always satisfied.

Examples of violations

- If the target treatment is only available to patients with certain socioeconomic characteristics (e.g., only wealthy patients due to insurance limitations) and these characteristics influence the outcome independently of the covariates in our model.
- If clinicians assign treatment based on patient-reported quality of life, that is not included in the model, but that does predict the outcome even after accounting for the variables that *are* included in the model.

Case study In the case study, treatment selection for protons follows a deterministic process based on baseline dysphagia, tumor location, planned dose under photons, and planned dose under protons, as described in section 2.3. If the model would have included all these variables used for selection, the condition would be satisfied. However, the model does not include the planned dose variables under protons, as proton plans were not available at the time model development data were collected, and arguably also do not seem intuitive predictors of outcomes under photon therapy. The question thus becomes: After conditioning on baseline dysphagia, tumor location, and planned photon dose variables, does knowledge of planned *proton dose* provide additional information about expected outcomes under photon therapy? Said differently, are there unmeasured variables that correlate with both the potential reduction in dose (from photons to protons) and dysphagia risk? The clearest example of such a variable would be the location of the tumor, which influences both the possible reduction in dose, as well as the risk of dysphagia. Tumor location is included in the model, although it may not be fully accounted for as it is only incorporated as broad categories. Another example is N-stage: patients with a higher N-stage often receive higher photon doses and may allow for greater potential dose reductions using protons. At the same time, a higher N-stage may be associated with a higher risk of dysphagia. Because N-stage is not included in the model, the predicted outcomes under photon therapy may be slightly *underestimated* in the proton-treated group, which could lead to an *underestimation* of the benefit of proton therapy.

To decide whether ignorability of treatment assignment holds, it may help to draw out assumptions about relevant causal relations. An example graphical representation of the case study is shown in Appendix D, which clarifies why the only risk for a violation of the condition is because the *proton* plan dose variables are not included in the model, whereas they did influence treatment assignment.

3.3.3 Condition 3: Consistency

This condition requires that for each patient receiving a treatment $T = t$, their observed outcome (Y_i) should match their potential outcome under that same treatment:

$$Y_i = Y_i(t) \quad \text{if } T_i = t$$

This reflects that treatment has a single well-defined meaning, i.e. all realistic ways of delivering a treatment are equivalent in terms of how they affect patient outcomes. In other words, once it is decided that a patient receives treatment $T = t$, *how* it is delivered (the specific equipment, clinician, or workflow used) should not change the patient’s outcome.

Why? Consistency allows us to connect what we observe in the data to the potential outcomes used in causal reasoning, such that we can interpret the estimated effect as the effect of treatment. Lack of consistency would result in an ambiguous treatment effect estimate, as you are effectively comparing different treatments under the same label.

Examples of violations

- If a patient is prescribed the same medication (e.g., insulin for diabetes), but the delivery method (e.g., insulin pump vs. manual injection) is not part of the treatment definition and would have led to different outcomes.
- If the clinical outcome of a patient depends strongly on the clinician’s skill or interpretation, even for the same treatment (e.g. a specific type of surgery, or radiotherapy planning).
- If an older and a newer radiation machine deliver the same therapy (e.g., photon radiotherapy) but approximate the planned dose with different levels of accuracy that would lead to different outcomes.

Case study In the case study, we assume that all meaningful variations of photon therapy are reflected by differences in dose distributions, which are included as covariates in the model. Some variations may still exist, even for the same planned dose distributions, e.g. due to variations in fractionation, but we believe it is a reasonable assumption that those variations would not result in a different grade of dysphagia at 6 months.

3.3.4 Condition 4: Positivity

This condition requires that all combinations of values of covariates in the model (both patient characteristics \mathbf{X} and treatment plan variables under the standard therapy $\mathbf{P}^{(0)}$), that could occur among the individuals in the post-introduction population that received the target treatment, could also occur in the pre-introduction population:

$$\text{Supp} \left[(\mathbf{X}, \mathbf{P}^{(0)}) \mid T = 1 \right]_{\text{post}} \subseteq \text{Supp} \left[(\mathbf{X}, \mathbf{P}^{(0)}) \right]_{\text{pre}}$$

Here, the support of a distribution refers to the set of values or combinations that the variables can take with non-zero probability. This is a somewhat theoretical condition: it does not require that all combinations actually occur in the data, but that there are no structural reasons preventing them from occurring.

Even if the positivity condition holds in theory, it may happen that some combinations are not observed in the available data simply due to chance, especially when the sample size is limited or certain subgroups are rare. This is known as a stochastic positivity violation [14, 15]. Although such stochastic violations do not invalidate the theoretical identifiability of the treatment effect (meaning that, with infinite data, the effect could still be estimated), they can cause similar practical challenges and should also be considered.

Why? To make reliable counterfactual predictions for patients in the target treatment group (i.e., those who received $T = 1$), their characteristics should be adequately represented in the population used to develop the model (the pre-introduction group). If certain combinations of patient features or treatment variables never appear in the pre-introduction population, the model needs to rely on extrapolation, essentially “guessing” outcomes for patients unlike anyone seen before. Depending on further parametric assumptions (e.g., linearity), this can lead to biased predictions or even make it impossible to make

predictions, for example, if a patient has a covariate value that was completely absent during model development.

Examples of violations

- If the pre-introduction population comes from a hospital that does not treat patients with oropharyngeal tumors, but the target treatment population includes such patients, and tumor location is included as a covariate in the model. In this case, the model cannot be used to make predictions in these patients.
- If the pre-introduction population is strictly adult (18+), but the target treatment population also contains teenagers, and age is used as a covariate in the model. Although this is an example of the positivity violation, it may not *necessarily* prevent making valid counterfactual predictions for the younger patients if the model used is capable of extrapolating to these younger ages (e.g., it can generate predictions for 16-year-olds). The validity of such extrapolations, however, depends heavily on whether the parametric assumptions underlying the model’s extrapolation are appropriate.

Case study In the case study, we expect this condition to hold because the pre-introduction sample comes from the same hospital as the sample of proton-treated individuals, differing only in time period. Examining the univariable distributions of covariates - baseline dysphagia, tumor location, photon plan dose parameters - we also observe substantial overlap (see Appendix A.2). Although looking at each variable independently does not tell the whole story, together with knowledge that the data comes from the same hospital, it does provide sufficient confidence that there are no issues with this condition.

3.3.5 Condition 5: Correct model specification

We require the outcome model \mathbf{m}_{pre} to be correctly specified (in terms of functional form) or sufficiently flexible to capture the conditional probability $P_{pre}(Y = 1 \mid \mathbf{X}, \mathbf{P}^{(0)})$, i.e., that there are model parameters β_0 such that $\mathbf{m}_{pre}(\mathbf{X}, \mathbf{P}^{(0)}; \beta_0) = P_{pre}(Y = 1 \mid \mathbf{X}, \mathbf{P}^{(0)})$. Note that this condition concerns the relationship between included covariates and the outcome, not whether additional covariates would improve the model.

Why? Since we use a statistical model to approximate conditional probabilities, misspecification of this model will lead to biased estimates of these probabilities and consequently biased ATT estimates.

Examples of violations

- If the age relates to the outcome in a non-linear way but is modeled with a linear term.
- If there is an interaction between tumor location and sex (e.g., a particular tumor location increases risk more in men) but the model is not flexible enough to capture this interaction. For example, it includes tumor location and sex as separate variables in a logistic regression without their interaction.

Case study The parametric model in the case study is relatively simple. While its functional form was chosen based on the absence of strong evidence for model misspecification in the pre-introduction data (i.e., no clear non-linearity in the association between dose parameters and the outcome was observed), it is still likely to be somewhat imperfectly specified. Since we are interested in the average predicted risk across the entire proton-treated population, some errors may average out, potentially limiting bias. Nonetheless, sensitivity analyses examining how different modeling choices affect the resulting ATT estimate would be advisable. We provide an example of such a sensitivity analysis in Appendix A.8.

3.4 Supportive evidence for validity

To estimate the average treatment effect among the treated, we make predictions about potential outcomes that are not observed and never will be. E.g., in the case study, for patients who received proton therapy, we estimated what would have happened if they had instead received photon therapy, which we never observe. This means that we can never directly validate our estimation strategy. However, in some settings, auxiliary data can be used for indirect validation.

For example, if some patients in the post-introduction population received the standard treatment, it is possible to assess how well the model predicts outcomes under the standard treatment in this sample. One

option is to use the same estimation procedure as before to estimate the ATT, but now in these individuals who actually received the standard treatment, i.e., estimate the average difference between observed and predicted outcomes under the standard treatment: $\frac{1}{N} \sum_{i: T_i=0, i \in \text{post}} (Y_i - \mathbf{m}_{\text{pre}}(\mathbf{X}_i, \mathbf{P}_i^{(0)}; \hat{\beta}))$.

If all assumptions hold, we expect this average difference to be close to zero. However, the converse is not necessarily true: good calibration in this group does not guarantee unbiased ATT estimates, since those are based on predictions for a different subgroup (the ones treated with protons), and this subgroup could systematically differ. Also, limited sample sizes will often limit the ability to draw strong conclusions. Still, a close agreement between observed and predicted outcomes provides supportive evidence for the robustness of the estimation strategy¹. If we do observe a systematic difference (e.g., consistent over- or underestimation), this would reflect violations of one or more conditions. One possibility is to interpret the difference as a *general* shift in outcome risk (e.g., due to setting or time), and use it to adjust the treatment effect estimates. But this again relies on assumptions that cannot be verified from the data alone.

Besides looking at patients from the post-introduction population who were treated with the standard treatment, more opportunities for obtaining supportive evidence for validity may exist. For example, in the case study, we can take advantage of the fact that both treatments (protons and photons) are forms of radiotherapy that mainly differ in the amount of delivered dose. If we assume that the relationship between dose and outcome, learned from photon therapy, also applies to proton therapy, we can assess model performance directly in the treated group (proton patients), even though the model was never directly trained to predict outcomes under proton therapy. We illustrate both types of validation in the case study and further discuss their interpretation in Appendix A.6 and Appendix A.7.

We described ways to gain evidence about the validity of the ATT estimation strategy, by comparing the average of observed and predicted outcomes of the model in auxiliary data. In prediction model literature, where validating models in an external data source is common practice, this is known as mean calibration or calibration-in-the-large [17]. Other methods for assessing calibration, such as calibration curves, can also be informative, provided the sample size is sufficient: if all underlying assumptions hold, we expect the model to be well-calibrated overall. In contrast, discrimination metrics such as the area under the ROC curve (AUROC) do, by themselves, not offer insight into the validity of assumptions of the validity of ATT estimation using the model. For example, a model may show poor discrimination (e.g., an AUROC near 0.5) due to a restricted or homogeneous case-mix, even if the assumptions are met.

4 Discussion

We have described the conditions under which a model fitted on pre-introduction patient data can be used to estimate the average treatment effect among the treated (ATT), by predicting counterfactual outcomes, i.e., what would have happened had the treated individuals received the standard treatment instead.

An important condition is that the population used to develop the model is sufficiently similar to the population on which the model is applied. Typically the target population, in which the model is applied is (1) a subset of patients eligible for the target treatment ('ignorability of treatment assignment'), and (2) from a different time and/or hospital ('transportability'). Whereas the populations are allowed to differ to some extent, the association between the covariates and the outcome is assumed to be stable. This is a strong assumption, and violations can arise in various ways. For example, the post-introduction population may have a higher prevalence of an unmeasured cause of the outcome, or treatment selection may depend on variables that are associated with the outcome, even after adjusting for covariates in the model. Another potential source of bias lies in model misspecification, e.g., when the model is not flexible enough to capture the true underlying predictor-outcome associations. In addition, issues may arise due to violations of consistency (if the treatment definition is ambiguous) or positivity, which could make it impossible to estimate the treatment effect.

Given these strong and largely untestable assumptions, why consider this approach at all? It is important to acknowledge that every study design relies on assumptions. In certain contexts, the conditions required for this approach may be more defensible than alternatives. For example, the most straight-

¹This validation strategy relates to the use of negative control populations [16], where the estimation procedure (in our case, the estimation of the ATT) is repeated on a group of individuals assumed to be similar to the population of interest, except that we expect the treatment effect to be zero.

forward analysis of the case study might be to compare the proportion of patients with dysphagia before and after the introduction of proton therapy. However, this strategy implicitly assumes that all changes in outcomes can be attributed to the introduction of proton therapy. In contrast, using a model for the outcome allows for adjustment for outcome shifts due to changes in covariate distributions (e.g., population aging) or changes in treatment plans for the standard treatment (e.g., improved dose planning). Another option could be to compare patients receiving the target versus the standard treatment during the post-introduction period, potentially using regression adjustment or weighting methods. However, these strategies require overlap in all relevant patient characteristics across groups, an assumption that is often violated in practice. For example, in our case study, treatment assignment follows a deterministic process based on expected treatment benefit, which will cause the proton-treated and photon-treated patients to have inherently different expected outcomes. By using a model trained on pre-introduction data, we can still indirectly compare with patients similar to the ones currently receiving the new treatment, but who received the standard treatment.

Besides being conceptually more applicable in some situations, as compared to other non-randomized approaches, the approach has other strengths. First, in principle, only the model needs to be shared or published; access to original patient-level data from the pre-introduction population is not required. However, this does raise challenges in estimating confidence intervals of the final ATT estimate, since the sampling variance in model development is often not readily available or easily propagated. Second, the approach is flexible with respect to the effect measure (e.g. risk ratio, odds ratio). Finally, a wide range of modeling strategies can be used, including flexible methods like neural networks when large pre-introduction datasets are available.

In the end, it is dependent on the context whether the approach is appropriate, and what level of evidence is needed to justify the untestable assumptions. Although we illustrated some ways to gather supportive evidence for the required assumptions, further work could explore the use of informed bias analyses and the derivation of plausible bounds on the causal effect [18–20]. This would help communicating not only sampling uncertainty but also uncertainty in the assumptions themselves, thereby improving transparency and practical relevance.

Furthermore, even when the approach is considered appropriate, there is still room for methodological refinement. Future research could examine strategies to improve model transportability, for example, using anchor regression or related methods [21, 22], or develop doubly robust variants of the estimator [23–25]. This could potentially lead to more robust and efficient estimation of the ATT.

Our work makes several contributions. First, although this strategy has been proposed and used as model-based clinical evaluation within radiotherapy research [4, 8–11], we are the first to formalize it as a causal inference method within the potential outcomes framework, and to formally derive and describe a set of sufficient conditions. Second, while similar strategies have been extensively discussed in causal inference literature [5–7, 26], the specific approach described here - using a model trained on pre-introduction data to estimate ATT in settings with personalized treatment plan information (such as dose planning in radiotherapy) - has not, to our knowledge, been formally described. By presenting this relatively specific approach within a formal causal inference framework, while illustrating it through a case study from its originating field, we aim to contribute both to the conceptual understanding and practical applicability of the approach. Finally, by complementing theory and examples with a detailed applied example in R we provide additional practical guidance.

There are several limitations to acknowledge. First, we did not cover the full range of model-based clinical evaluation strategies used in radiotherapy. For example, Rwigema et al. [10] used a model trained on a mixed cohort of patients treated with multiple radiation technologies (instead of only the standard technology), assuming a shared dose-response relationship. We chose not to focus on this variant because it is highly specific to radiotherapy, and our goal was to present a more generalizable approach that could be useful beyond this field. However, the core assumptions and conditions we described are largely the same or can be easily adapted. Furthermore, although we aimed to maintain practical applicability, our primary focus was causal identifiability. To concentrate on conditions for causal identifiability, the case study, while inspired by the Dutch context, represented a simplified version of reality. We acknowledge that real-world analyses often involve additional complexities, such as missing data, competing risks, measurement error, or limited sample sizes, which may introduce further challenges.

To conclude, we have discussed the theory behind an estimation approach for the ATT: using a model trained on pre-introduction data, to predict counterfactual outcomes in patients treated with the target treatment. While this approach could offer practical and conceptual advantages, it is only valid under specific conditions. When using the approach, we recommend a systematic evaluation of all these required conditions in collaboration with clinical and methodological experts. Where possible, empirical evidence should be gathered to determine whether the conditions are likely to hold, and if not, what

impact this has on the conclusions. These practices will improve the credibility and transparency of the treatment effect estimates and ultimately support more informed clinical decisions.

Declarations

Funding This project has received funding from ZonMw HTA Methodology grant number 10580012210025 (WhyMBA). This dissemination reflects only the authors' views and the Commission is not responsible for any use that may be made of the information it contains.

Author contribution ES and AML conceptualized the initial study. LMM drafted the manuscript and created the figures and code under supervision of ES, AML and KGM. All other authors critically reviewed the manuscript and contributed to revisions. ES is the guarantor of this work. The corresponding author (LMM) attests that all listed authors meet authorship criteria and that no others meeting the criteria have been omitted

Conflict of interest All authors have completed the ICMJE uniform disclosure form at <http://www.icmje.org/disclosure-of-interest/> and declare: support from ZonMW for the submitted work. RN reports research grants and honoraria from Elekta, Varian, Accuray, Sensius, dr Sennewald, MSD, and GSK, paid to their institution; JAL reports research grants from the European Union and the Dutch Cancer Society, consulting fees and honoraria paid to UMCG Research BV by IBA, is chair of the safety monitoring board of the UPGRADE trial (UMC Nijmegen), member of scientific advisory committees for IBA and RaySearch, and reports departmental research collaborations with IBA, RaySearch, Elekta, Mirada, and Siemens; Other authors declare no financial relationships with any organisations that might have an interest in the submitted work in the previous three years; no other relationships or activities that could appear to have influenced the submitted work.

Ethics approval and consent to participate This study did not involve real patient data; therefore, ethical approval was not required.

References

- [1] Eduardo Hariton and Joseph J. Locascio. "Randomised controlled trials—the gold standard for effectiveness research". In: *BJOG : an international journal of obstetrics and gynaecology* 125.13 (Dec. 2018), p. 1716. ISSN: 1470-0328. DOI: 10.1111/1471-0528.15199.
- [2] Robbe Saesen et al. "Defining the role of real-world data in cancer clinical research: The position of the European Organisation for Research and Treatment of Cancer". In: *European Journal of Cancer* 186 (June 2023), pp. 52–61. ISSN: 0959-8049. DOI: 10.1016/j.ejca.2023.03.013.
- [3] *Evaluation of new technology in health care: in need for guidance for relevant evidence*. Tech. rep. Amsterdam: KNAW, 2014. URL: <https://www.knaw.nl/publicaties/evaluation-new-technology-health-care>.
- [4] Johannes A. Langendijk et al. "Clinical Trial Strategies to Compare Protons With Photons". In: *Seminars in Radiation Oncology*. Proton Radiation Therapy 28.2 (Apr. 2018), pp. 79–87. ISSN: 1053-4296. DOI: 10.1016/j.semradonc.2017.11.008.
- [5] MA Hernán and JM Robins. *Causal Inference: What if*. Boca Raton: Chapman & Hall/CRC, 2020. ISBN: 1-4200-7616-7.
- [6] Issa J. Dahabreh et al. "Generalizing causal inferences from individuals in randomized trials to all trial-eligible individuals". eng. In: *Biometrics* 75.2 (June 2019), pp. 685–694. ISSN: 1541-0420. DOI: 10.1111/biom.13009.
- [7] Timothy L Lash et al. *Modern Epidemiology*. 4th ed. Wolters Kluwer, 2021. ISBN: 978-1-4511-9328-2.

- [8] Tineke W.H. Meijer, Dan Scandurra, and Johannes A. Langendijk. “Reduced radiation-induced toxicity by using proton therapy for the treatment of oropharyngeal cancer”. In: *British Journal of Radiology* 93.1107 (Mar. 2020), p. 20190955. ISSN: 0007-1285. DOI: 10.1259/bjr.20190955.
- [9] Nataniel H. Lester-Coll and Danielle N. Margalit. “Modeling the Potential Benefits of Proton Therapy for Patients With Oropharyngeal Head and Neck Cancer”. In: *International Journal of Radiation Oncology*Biophysics* 104.3 (July 2019), pp. 563–566. ISSN: 0360-3016. DOI: 10.1016/j.ijrobp.2019.03.040.
- [10] Jean-Claude M. Rwigema et al. “A Model-Based Approach to Predict Short-Term Toxicity Benefits With Proton Therapy for Oropharyngeal Cancer”. In: *International Journal of Radiation Oncology*Biophysics* 104.3 (July 2019), pp. 553–562. ISSN: 0360-3016. DOI: 10.1016/j.ijrobp.2018.12.055.
- [11] Miranda E. M. C. Christianen et al. “Swallowing sparing intensity modulated radiotherapy (SW-IMRT) in head and neck cancer: Clinical validation according to the model-based approach”. eng. In: *Radiotherapy and Oncology: Journal of the European Society for Therapeutic Radiology and Oncology* 118.2 (Feb. 2016), pp. 298–303. ISSN: 1879-0887. DOI: 10.1016/j.radonc.2015.11.009.
- [12] Johannes A. Langendijk et al. “Selection of patients for radiotherapy with protons aiming at reduction of side effects: The model-based approach”. In: *Radiotherapy and Oncology* 107.3 (June 2013), pp. 267–273. ISSN: 0167-8140. DOI: 10.1016/j.radonc.2013.05.007.
- [13] Donald B. Rubin. “Causal Inference Using Potential Outcomes: Design, Modeling, Decisions”. In: *Journal of the American Statistical Association* 100.469 (2005). Publisher: [American Statistical Association, Taylor & Francis, Ltd.], pp. 322–331. ISSN: 0162-1459. URL: <https://www.jstor.org/stable/27590541>.
- [14] Maya L Petersen et al. “Diagnosing and responding to violations in the positivity assumption”. In: *Statistical methods in medical research* 21.1 (Feb. 2012), pp. 31–54. ISSN: 0962-2802. DOI: 10.1177/0962280210386207.
- [15] Paul N. Zivich, Stephen R. Cole, and Daniel Westreich. *Positivity: Identifiability and Estimability*. arXiv:2207.05010 [stat]. July 2022. DOI: 10.48550/arXiv.2207.05010.
- [16] Marco Piccininni and Mats Julius Stensrud. “Using Negative Control Populations to Assess Unmeasured Confounding and Direct Effects”. en-US. In: *Epidemiology* 35.3 (May 2024), p. 313. ISSN: 1044-3983. DOI: 10.1097/EDE.0000000000001724.
- [17] Ben Van Calster et al. “Calibration: the Achilles heel of predictive analytics”. In: *BMC Medicine* 17.1 (Dec. 2019), p. 230. ISSN: 1741-7015. DOI: 10.1186/s12916-019-1466-7.
- [18] “Bias Analysis”. In: *Modern Epidemiology*. 4th ed. Wolters Kluwer, 2021, pp. 1556–1616. ISBN: 978-1-4511-9328-2.
- [19] Trang Quynh Nguyen et al. “Sensitivity analysis for an unobserved moderator in RCT-to-target-population generalization of treatment effects”. In: *The Annals of Applied Statistics* 11.1 (Mar. 2017). Publisher: Institute of Mathematical Statistics, pp. 225–247. ISSN: 1932-6157, 1941-7330. DOI: 10.1214/16-AOAS1001.
- [20] Issa J. Dahabreh et al. “Sensitivity analysis using bias functions for studies extending inferences from a randomized trial to a target population”. en. In: *Statistics in Medicine* 42.13 (2023), pp. 2029–2043. ISSN: 1097-0258. DOI: 10.1002/sim.9550.
- [21] Xinwei Shen, Peter Bühlmann, and Armeen Taeb. *Causality-oriented robustness: exploiting general additive interventions*. arXiv:2307.10299 [stat]. July 2023. DOI: 10.48550/arXiv.2307.10299.
- [22] Dominik Rothenhäusler et al. “Anchor Regression: Heterogeneous Data Meet Causality”. In: *Journal of the Royal Statistical Society Series B: Statistical Methodology* 83.2 (Apr. 2021), pp. 215–246. ISSN: 1369-7412. DOI: 10.1111/rssb.12398.

- [23] Heejung Bang and James M. Robins. “Doubly Robust Estimation in Missing Data and Causal Inference Models”. In: *Biometrics* 61.4 (Dec. 2005), pp. 962–973. ISSN: 0006-341X. DOI: 10.1111/j.1541-0420.2005.00377.x.
- [24] Mark J. Van Der Laan and Sherri Rose. *Targeted Learning: Causal Inference for Observational and Experimental Data*. en. Springer Series in Statistics. New York, NY: Springer, 2011. ISBN: 978-1-4419-9781-4 978-1-4419-9782-1. DOI: 10.1007/978-1-4419-9782-1.
- [25] James M. Robins, Andrea Rotnitzky, and Lue Ping Zhao. “Estimation of Regression Coefficients When Some Regressors Are Not Always Observed”. In: *Journal of the American Statistical Association* 89.427 (1994). Publisher: [American Statistical Association, Taylor & Francis, Ltd.], pp. 846–866. ISSN: 0162-1459. DOI: 10.2307/2290910.
- [26] Stuart J. Pocock. “The combination of randomized and historical controls in clinical trials”. In: *Journal of Chronic Diseases* 29.3 (Mar. 1976), pp. 175–188. ISSN: 0021-9681. DOI: 10.1016/0021-9681(76)90044-8.