

# Estimation of dose–response functions for longitudinal data using the generalised propensity score

Erica EM Moodie<sup>1</sup> and David A Stephens<sup>2</sup>

Statistical Methods in Medical Research  
21(2) 149–166

© The Author(s) 2010

Reprints and permissions:

sagepub.co.uk/journalsPermissions.nav

DOI: 10.1177/0962280209340213

smm.sagepub.com



## Abstract

In a longitudinal study of dose–response, it is often necessary to adjust for confounding or non-compliance, which may otherwise compromise the estimation of the true effect of a treatment. Using an approach based on the generalised propensity score (GPS) – a generalisation of the classical, binary treatment propensity score – it is possible to construct a balancing score that provides an estimation procedure for the true (unconfounded) direct effect of dose on response. Previously, the GPS has been applied only in a single interval setting; in this article, we extend the GPS methodology to the longitudinal setting to estimate the direct effect of a continuous dose on a longitudinal response. The methodology is applied to two simulated examples, and a real longitudinal dose–response investigation, the Monitored Occlusion Treatment of Amblyopia Study (MOTAS). In the treatment of childhood amblyopia, a common ophthalmological condition, occlusion therapy (patching) was for many decades the standard medical treatment, despite the fact that its efficacy was not quantified. MOTAS was revolutionary, as it was the first study to obtain precise measurements of the amount of occlusion each study participant received over the course of the study.

## Keywords

causal inference, confounding, continuous dose-response, generalized propensity score, longitudinal data, monitored occlusion treatment of amblyopia study, noncompliance

## 1 Introduction

In observational studies of the efficacy of a treatment, there is the potential for bias in the estimation of the treatment effect whenever the treatment dose level received is influenced by subject-specific covariates. Randomised trials, particularly those where treatment is given over time in several treatment intervals, must also contend with partial or total non-compliance, which arguably renders the trial an observational study of the effect of received treatment (though still a

<sup>1</sup>Department of Epidemiology & Biostatistics, McGill University, 1020 Pine Ave W., Montreal, QC, Canada.

<sup>2</sup>Department of Mathematics and Statistics, McGill University, 805 Sherbrooke Str W., Montreal, QC, Canada.

### Corresponding author:

Erica EM Moodie, Department of Epidemiology & Biostatistics, McGill University, 1020 Pine Ave W, Montreal, QC H3A 1A2 Canada  
Email: erica.moodie@mcgill.ca

randomised study of assigned treatment). Statistical analyses in the face of non-compliance have often relied on intention-to-treat or as-treated analyses, which respectively ignore the dose actually received or do not account for the informative nature of non-compliance. The aim of this article is to provide a framework for examining the direct effect of treatment given *over time* with either incomplete adherence to prescribed dose or at a patient-controlled level.

### 1.1 Longitudinal observational dose–response studies

There are two primary effects of interest in longitudinal data where an exposure such as a dose of a treatment is received over time: the *direct effect* the treatment received at one instance on a response some time in the future, or the *total* or cumulative (direct and indirect, i.e. mediated) effect of the treatment. In both cases, the data are collected over time, with each participant providing not only baseline data but also (potentially time-varying) covariate, treatment, and response data over a number of observation intervals. If interest lies in the direct effect of a treatment dose on the next-measured response, examining the by-interval changes (i.e. taking as the response, the vector of differences in, say, health outcomes between successive measurements) may simplify the analysis by reducing or even removing the serial correlation which may be present in the data. Using the vector of differences in response also has the advantage of allowing the analyst to fit a common curve in all intervals (e.g. a linear dose–response with a common intercept as well as slope) provided the treatment effect is not modified by time. A recent review of the main statistical issues that arise when estimating causal effects from observational longitudinal data can be found in Arjas and Parner.<sup>1</sup>

In this article, we focus on the by-interval changes, and estimate the direct effect of dosing that may be attributed to a unit of dose in an interval. To account for subject-controlled treatment level, which we interpret as *non-* or *partial* compliance, a potential source of bias in the estimation of treatment effect, we develop methods that address issues of confounding and non-compliance using a *balancing score* approach based on the generalised propensity score (GPS)<sup>2,3</sup> for a continuous treatment that controls for sources of such bias. The GPS has received relatively little attention in statistical circles; however, see Flores<sup>4</sup> for an economics application.

As discussed in Rubin,<sup>5</sup> observational studies and randomised trials are part of a continuum, where confounding of the effect of dose received in randomised studies may arise due to non-compliance. In the context of randomised trials, causal methods have been proposed which are based on principal stratification by the compliance score;<sup>6–8</sup> see Joffe et al.<sup>9</sup> for a comprehensive explanation and discussion. In such an approach, the compliance score, a predictive model for compliance with assigned treatment given baseline covariates, is used to estimate the complier average causal effect, that is, the effect of treatment among individuals who would comply with a prescribed treatment dose. Adjustments for compliance based only on pre-treatment variables will typically be unable to identify the effect of a treatment taken over time, as compliance may vary as a function of response to treatment.

### 1.2 Causal methods for repeated measures data

Causal methods for estimating treatment effects on a univariate end-of-study outcome (e.g. depression score 4 years after diagnosis) from repeated measures data are available but not readily implemented for continuous doses. For example, a G-computation procedure was proposed in Neugebauer and van der Laan,<sup>10</sup> which was computationally expensive, and so the authors discretised the continuous treatment. Another causal procedure for repeated measures

data uses marginal structural models (MSMs).<sup>11,12</sup> MSMs use inverse weighting using the treatment mechanism model to estimate the marginal effect of time-varying treatment regimes on a univariate response. MSMs were originally designed to estimate the effects of static (not covariate-adapted) time-varying treatments, but were recently proposed as a method to compare dynamic regimes which may have been found by other estimation procedures<sup>13</sup> and further extended to estimate the optimal dynamic regime.<sup>14–16</sup> MSMs can be used when treatments are measured on a continuous scale, although this is not common in practice. Both the G-computation and the MSM approach are designed to examine total effects of dosing regimes on end-of-study outcomes, rather than the direct effects that are addressed in this article.

### 1.3 Objectives and structure of the article

In this article, we extend a causal modelling approach to account for the within-person correlation of responses, the within-person correlation of doses, and the potential confounding of the direct effect of dose on response by previous doses and responses. The article is structured as follows: Sections 2.1–2.4 introduce notation and the GPS methodology, the quantity which is the focus of interest, the average potential outcome, and an algorithm to estimate this quantity. Section 2.5 extends the balancing score approach to incorporate the complexities of a repeated measures structure with time-varying covariates. Section 3 contains two simulated examples illustrating the methodology where some analytical calculations are possible. In section 4, the first causal analysis of the Monitored Occlusion Treatment of Amblyopia Study (MOTAS) is undertaken, where the dose–response relationship between occlusion and improvement in visual acuity is quantified. We conclude with a discussion in section 5.

## 2 A balancing score approach to estimating a dose–response relationship

To ascertain the true direct effect of dose, a causal analysis, which accounts for the possible confounding of treatment or dose effect by other measured covariates, is necessary. One tool used to account for possible confounding relationships between occlusion treatment and other covariates is the GPS,<sup>2,3,17</sup> a constructed variable that can be used (in a regression analysis or via stratification) for removal of bias in the estimation of the treatment effect. A regression, which includes the GPS, does not directly provide a parameter that may be interpreted causally; however, it can be used to obtain estimates of the potential response to a dose, which do have a causal interpretation.

### 2.1 Notation

Suppose that we have collected data repeatedly on  $N$  individuals, so that  $n_i$ ,  $i=1, \dots, N$  measurements are available for each subject. We denote the total number of data points by  $n = \sum_{i=1}^N n_i$ , the response for individual  $i$  at time  $j$  by  $Y_{ij}$ , the treatment dose by  $D_{ij}$ , and other possibly confounding covariates by  $X_{ij}$ ; we denote observed values of these random variables  $y_{ij}$ ,  $d_{ij}$  and  $x_{ij}$ , respectively. We define  $\mathcal{D}$ , a bounded interval in  $\mathbb{R}$ , to be the set of possible doses.

Initially, for simplicity, we restrict attention to the case where  $n_i = 1$ , and drop the interval-specific subscript  $j$ , so the response data are simply  $Y_i \in \{Y_{ij} : i=1, \dots, N, j=1\}$  and  $n = N$ . We address the general longitudinal case in section 2.5.

## 2.2 The generalised propensity score

The causal analysis is formulated through the use of *potential* or *counterfactual* outcomes. A potential outcome is a value of the response that would result if a subject were to receive a specified treatment dose, not necessarily the same dose that they received in the study. We denote by  $Y_i(d) \equiv Y_i(D_i = d)$  the potential response random variable resulting from a dose  $D_i = d$  taken in an interval, and write  $y_i(d)$  for the observed version. Throughout the article,  $d$  without a subscript will indicate a potential dose. Potential outcomes adhere to the *axiom of consistency*: the actual and potential response are equal when the regime in question is the dose actually received, that is,  $y_i(d) = y_i$  if  $d_i = d$ .

As with all models for observational data, causal models require certain modelling assumptions to be appropriately specified.<sup>11,18</sup> Specifically, we make the *stable unit treatment value assumption*,<sup>19</sup> which states that a subject's outcome is not influenced by other subjects' treatment allocation. We further assume *weak unconfoundedness*: for all  $d \in \mathcal{D}$ , the potential outcome  $Y_i(d)$  and the dose received  $D_i$  are presumed conditionally independent given the covariates  $X_i$ , that is  $Y_i(d) \perp D_i | X_i$ . Informally, weak unconfoundedness implies that the mechanism which dictates response to any specific (potential) dose  $d$  and the mechanism by which dose is allocated, are probabilistically independent, conditional on the covariates. See Hirano and Imbens<sup>3</sup> for formal arguments. Here we note that this assumption is no stronger than that required for unbiased, causal inference from a standard regression approach: that is, in any regression, all confounding variables must have been measured and included in the response model to ensure unbiased estimation of the true effect of an exposure on the response. Thus, the GPS methodology requires the same assumption as standard covariate adjustment.

Following Imbens<sup>2</sup> and Hirano and Imbens,<sup>3</sup> we define the GPS,  $r_i = r(d, x_i)$  for any dose  $d$  and observed covariate values  $x_i$  by

$$r_i = r(d, x_i) = f_{D_i|X_i}(d|x_i), \quad (1)$$

that is, the conditional density function for  $D_i$  given  $X_i = x_i$  evaluated at  $D_i = d$ ; the random variable  $R_i = r(d, X_i)$  denotes a corresponding random quantity for fixed  $d$ . Note that these are potential quantities that may be evaluated at  $d = d_i$  and  $x = x_i$  to yield the *observed* GPS; we reserve the notation  $\hat{r}_i$  for this special case, that is we define

$$\hat{r}_i = r(d_i, x_i) = f_{D_i|X_i}(d_i|x_i).$$

The GPS is an extension of the propensity score<sup>20</sup> to continuous treatments. In this article, we regard the construction of the conditional density as a regression problem, and regress  $D_i$  on  $X_i$ ; that is, we fit a regression model to the pairs  $(x_i, d_i)$ ,  $i = 1, \dots, N$  in order to be able to compute  $r(d, x_i)$  for any  $d$ . Note that, for two doses  $d^1 \neq d^0$  we will have

$$r(d^1, x_1) \neq r(d^0, x_0),$$

in general, even if  $x_1 = x_0$ .

The GPS quantities,  $R_i$  and  $r_i$ , form part of the bias removal strategy. As is shown by Hirano and Imbens,<sup>3</sup> the GPS random quantity  $R_i$  has two properties that render it useful in causal inference problems. First,  $R_i$  acts as a *balancing* score, in that  $D_i$  and  $X_i$  are conditionally independent given  $R_i$ ; in particular, within strata of  $R_i$ , the distribution of dose is (approximately) the same irrespective of the value of the covariate. Secondly, for any  $d$ , the

distribution of the treatment dose is conditionally independent of the potential response, given the propensity score,

$$Y_i(d) \perp D_i | r(d, X_i)$$

that is, we have weak unconfoundedness of  $Y_i(d)$  and  $D_i$  given  $R_i$ .

The first point, that  $R_i$  breaks the dependence between  $D_i$  and  $X_i$ , is the crucial factor that permits causal inference; the second point permits simplified modelling. We shall see that with natural extensions, both features carry over to the longitudinal setting.

### 2.3 Average potential outcomes

The causal effect in a single-interval study on which we focus is the marginal effect of dose on the response. Specifically, a typical quantity of interest in causal dose–response modelling is the average potential outcome (APO) at dose level  $d$ ,  $\mu(d) = E[Y_i(d)]$ , which traces the causal dose–response relationship as  $d$  varies in  $\mathcal{D}$ . That is, we are interested in  $\mu(d)$  for  $d \in \mathcal{D}$  where the response  $Y_i$  may be a health rating or a change in health score from baseline. We make the assumption that  $\mu_i(d) = \mu(d)$  for all  $i$ , so that

$$\mu(d) = E[Y_i(d)] = E_{X_i}[E[Y_i(d)|X_i]].$$

To report the causal effect of interest, we first examine the *conditional* average causal effect of dose, defined as the difference in expected outcomes for two potential dose levels  $d^0, d^1$  for fixed covariate values  $X_i = x$ , that is

$$E[Y_i(d^1)|X_i = x] - E[Y_i(d^0)|X_i = x]. \quad (2)$$

The *marginal* average causal effect is the expectation of this quantity over the distribution of different  $X$  values in the study population, provided all confounding covariates are included in  $X$ . Such an approach provides interpretable estimates of the causal dose–response relationship, both conditionally and after marginalising over the distribution of  $X$ . That is, marginalising Equation (2) with respect to the distribution of  $X$  yields an estimate of  $\mu(d^1) - \mu(d^0)$ . However, note that to estimate the marginal average effect, it is necessary first to compute the conditional expectation of response given dose and covariates, and then to compute the conditional expectation of the covariates given the dose. This latter step requires knowledge of the conditional expectation of  $X$  at each  $d \in \mathcal{D}$ .

The modelling of outcome on dose and the GPS,  $R_i$  (rather than  $X_i$ ), returns an estimate of

$$E[Y_i(d^1)|R_i = r] - E[Y_i(d^0)|R_i = r], \quad (3)$$

and its population average. Thus GPS does not return estimates of the (causal) quantity in Equation (2), but does yield a bias-removal strategy: we examine the conditional distribution of  $Y_i(d)$  for fixed  $d$  given  $R_i = r(d, X_i)$ , rather than the conditional distribution given for fixed  $d$  given  $X_i$ , and recover a consistent estimator of the dose–response relationship by averaging over the empirical distribution of  $R_i$  in the population.

More specifically, Equation (3) facilitates consistent estimation of  $\mu(d)$ , as we may average the conditional expectations over the distribution of  $R_i$  if the balancing property holds, that is, if within

strata of  $R_i$ , the conditional distribution of  $D_i$  does not depend on  $X_i$ . The adequacy of any proposed propensity score model rests on whether or not balance is achieved, but this can be checked by standard exploratory statistical methods. In particular, strata of  $R_i$  may be formed; with the strata, the dose–response relationship may be estimated by, for example, regression; the relationship between doses and responses may then be averaged over the strata.

## 2.4 An algorithm for estimating the APO

The role of the propensity score in estimating the APO is made clear by the identity given in Imbens.<sup>2</sup> Assuming that  $\mu_i(d) = \mu(d)$  for all  $i$ , we may drop the dependence on  $i$  and write  $\mu(d) = E[Y(d)]$ . This marginal expectation can be achieved via the iterated expectation  $E_X[E[Y(d)|X]]$ ; however, this requires a model for the response as a function of the potentially very high dimensional covariate vector,  $X$ . Instead, we note that

$$E_X[E[Y(d)|X]] \equiv E_X[E[Y(d)|X, r(d, X)]] = E_X[E[Y(d)|r(d, X)]],$$

so that the iterated expectation over  $Y$  given  $X$ , then  $X$ , that is computed at the fixed dose  $d$ , is replaced by an iterated expectation over  $Y$  given  $R$ , then  $R$  (at the fixed dose  $d$ ), utilising the fact that for fixed  $d$  and  $X$ ,  $R$  is completely determined. Because of the properties of the GPS, both the internal and external expectations

$$E[Y(d)|r(d, X)] \quad \text{and} \quad \{E_X[E[Y(d)|r(d, X)]]\},$$

can be estimated consistently and without bias from the sample data for any potential  $d$  using the corresponding empirical averages; we require in turn models for  $Y$  given  $R$ , and then  $D$  given  $X$  to get  $R$ . In contrast, for the ‘direct’ approach, we have the two expectations

$$E[Y(d)|X] \quad \text{and} \quad E_X[E[Y(d)|X]]$$

and require the correct model for  $Y$  given  $D$  and  $X$ , and then the correct model for  $X$  given  $D = d$ , at all  $d \in \mathcal{D}$ . The latter of these models is considerably more difficult to learn from observed data, especially when  $X$  is high dimensional. Thus, in general, the benefit afforded by the GPS approach lies in the simplification of the model specification for the response: in the GPS approach, response is modelled as a function of dose and a one-dimensional summary of the covariates, i.e. the GPS, in contrast to standard regression modelling where response must be specified as a function of dose and the possibly very high dimensional covariate vector. We note, however, that it is possible and may in some cases desirable to model response not only as a function of dose and the GPS, but also some small subset of the covariates  $X$ .

We outline the estimation procedure of Hirano and Imbens<sup>3</sup> tailored specifically to the context discussed in this article. We presume parametric models for the two components of the procedure:

- (I) **Form the GPS model:** Construct a model,  $f_{D_i|X_i}(d|x_i, \beta)$ , for dose  $D_i$  given covariates  $X_i$  by regressing the observed doses  $d_i$  on the observed covariates  $x_i$ . Estimate parameters  $\beta$ , from the observed dose and covariate data  $\{(d_i, x_i), i = 1, \dots, N\}$ .
- (II) **Compute the fitted GPS values:** Compute the estimated GPS,  $\hat{r}_i = f_{D_i|X_i}(d_i|x_i, \hat{\beta})$ .
- (III) **Form the observable model:** Construct a model  $f_{Y(d)|X, R}(y(d)|x, r, \alpha)$  for  $Y_i$  for given dose value  $d$ , covariates  $X_i$  and propensity score  $R_i$ , using a regression approach and a model with



parameters  $\alpha$ . Note that we would only usually include a subset of the  $X_i$  in this model. Estimate parameters  $\alpha$  using the observed data and estimated GPS values  $\{(y_i, d_i, x_i, \hat{r}_i), i = 1, \dots, N\}$ .

(IV) **Estimate the APO:** For  $d \in \mathcal{D}$ , estimate the APO at dose level  $d$  by

$$\hat{\mu}(d) = \hat{E}[Y_i(d)] = \frac{1}{N} \sum_{i=1}^N E[Y_i(d) | X_i = x_i, r_i = r(d, x_i), \hat{\alpha}],$$

where  $r_i$  is evaluated from the model in (I) at  $\beta = \hat{\beta}$ .

The result of this procedure,  $\hat{\mu}(d), d \in \mathcal{D}$ , is the GPS-adjusted estimated dose–response function. Uncertainty bounds can be obtained analytically from the parametric analysis, or by a bootstrap procedure where individuals are sampled with replacement, which may be more straightforward in the longitudinal setting. Justification for the GPS APO estimation procedure is given in Hirano and Imbens,<sup>3</sup> and is extended to the repeated measures case in section 2.5. The two key conditional models  $f_{D|X}(d|x, \beta)$  and  $f_{Y(d)|X, R}(y(d)|x, r, \alpha)$  or the corresponding conditional moments, must be user-specified, but the adequacy of both components can be assessed in a straightforward statistical fashion. Note that we have allowed for the possibility that the observable model is formulated including terms in  $X$ , as discussed in Imai and Van Dyk,<sup>17</sup> but in many settings (including the analysis of section 4), the observable model will be formulated in terms of the dose and the GPS only.

We note that any one-to-one function of the GPS provides the desired balancing property; in addition, categories defined by discretising  $R_i$  may also provide sufficient balance to remove most of the bias due to confounding variables. In particular, an alternative approach proposed by Imai and Van Dyk<sup>17</sup> suggests that the APO may be approximated by estimating the dose–response effect within strata defined by the linear predictor of the treatment density function, and then combining these estimates to form a single, weighted average. This approach is straightforward to implement and often provides an estimate of the dose–response relationship that has little or no residual bias, although it may be less efficient than the regression approach described above.

## 2.5 The GPS for repeated measures data

In the case of dose–response estimation from repeated measures or multi-interval data, the potential patterns of confounding are more complex than that can be dealt with using a univariate GPS approach. In this section, we formulate a GPS approach suitable for the analysis of repeated measures response data with interval-dependent dosing. In the repeated measures setting, we no longer wish to ignore the correlation structure in the data, and so we return to the use of notation that makes this explicit. We therefore have that  $Y_{ij}, i = 1, \dots, N, j = 1, \dots, n_i$  is the response for individual  $i$  in interval  $j$ ; dose and covariate variables are similarly subscripted.

We assume that the marginal distribution of the counterfactual response,  $Y_{ij}(d)$ , is not modified by time, so that we have a single dose–response function  $\mu(d) = E[Y_{ij}(d)] = E[Y(d)]$  to estimate. The method is amenable to more general models, however, including those in which there are treatment-by-time interactions.

Repeated measures, or multivariate, data require a modification of the GPS procedure to account for confounding of the direct effect of dose  $D_{ij}$  on response  $Y_{ij}$  by previous doses and responses to treatment. Thus, we allow covariates  $X_{ij}$  at time  $j$  to include treatment doses and responses to treatments for person  $i$  at times  $1, 2, \dots, j-1$ . We denote the history of covariates, response, and

previous doses by  $\check{X}_{ij} = (X_{1j}, \dots, X_{ij})^T$ , and let  $R_{ij} = r(d, \check{X}_{ij})$ . Furthermore, we modify the notion of weak unconfoundedness to what we term *sequential* weak unconfoundedness; we assume

$$Y_{ij}(d) \perp D_{ij} | \check{X}_{ij}.$$

That is, at each interval, assignment to dose  $D_{ij}$  is weakly unconfounded with the response during interval  $j$  given covariates, previous response and dose values measured up to the start of the  $j$ -th interval.

We now demonstrate that the multivariate GPS (MGPS) procedure – the GPS constructed to allow for multivariate or repeated doses – retains the desired balancing properties of the univariate approach in a repeated measures setting. The results for the single interval setting can be recovered from the theorems as special cases.

**Theorem 1 (Weak unconfoundedness given the MGPS):** *Suppose that assignment to treatment in the  $j$ -th interval is sequentially weakly unconfounded given variables  $\check{X}_{ij}$  that occurred prior to treatment in the current interval (and may include previous treatment doses). Then, for every dose  $d$ ,*

$$Y_{ij}(d) \perp D_{ij} | R_{ij},$$

*that is, for  $d \in \mathcal{D}$ , current potential response  $Y_{ij}(d)$  is conditionally independent of the distribution of dose received  $D_{ij}$  given the MGPS  $R_{ij}$ , for all  $i$  and  $j$ .*

**Theorem 2 (Bias removal of the MGPS procedure):** *Suppose that  $\mu(d) = E[Y_{ij}(d)] = E[Y(d)]$  is the marginal mean of interest. For interval  $j$ , consider the mean*

$$\beta(d, r) = E[Y_{ij}(d) | R_{ij} = r(d, \check{X}_{ij}) = r]$$

*that conditions on the MGPS. The APO, obtained by averaging  $\beta(d, r)$  over the observed distribution of the covariates  $\check{X}_{ij}$ , is an unbiased estimator of the dose–response function  $\mu(d)$ .*

Proofs of these theorems follow in a straightforward fashion from the results in Hirano and Imbens,<sup>3</sup> and are included for completeness in the Appendix.

By Theorem 2, applying the bias removal result sequentially to each interval, we obtain an unbiased estimator of  $\mu(d)$  after pooling results over all intervals, by taking the expectation in turn over  $\check{X}_{i1}, \check{X}_{i2}, \dots$ . Note that a ‘univariate’ GPS analysis that does not construct a GPS by conditioning on  $\check{X}_{ij} = \check{X}_{ij}$  for each  $j$  may not achieve bias removal.

We have carried out extensive testing of the MGPS approach and performed comparisons with non-causal and standard GPS methods. Our examples demonstrate the importance of the use of the multivariate extension of the GPS provided in this article.

### 3 Simulation studies

#### 3.1 Simulation I: Non-linear, non-additive treatment effect

We extended the artificial example of Hirano and Imbens<sup>3</sup> to a two-interval, two-confounder setting. In this section, we drop dependence on subject index  $i$  for convenience; we have covariates measured in two intervals, the first subscript will denote interval and the second subscript will denote variable.

**Data generation:** Suppose that at the first and second interval, we have



$$Y_1(d)|X_{11}, X_{12} \sim \mathcal{N}(d + (X_{11} + X_{12}) \exp[-d(X_{11} + X_{12})], 1)$$

$$Y_2(d)|X_{21}, X_{12} \sim \mathcal{N}(d + (X_{21} + X_{12}) \exp[-d(X_{21} + X_{12})], 1),$$

and that the marginal distributions of each of  $X_{11}$ ,  $X_{12}$  and  $X_{21}$  are all unit exponential. Let  $D_1 \sim \exp(X_{11} + X_{12})$ ,  $D_2 \sim \exp(X_{21} + X_{12})$ . The marginal mean of the response in either interval is identical. As in Hirano and Imbens,<sup>3</sup> the APO can be obtained by integrating out the covariates analytically, yielding

$$\mu(d) = d + \frac{2}{(1+d)^3}.$$

A multivariate GPS (MGPS) analysis will involve the concatenated vector  $R^M = (R_1, R_2)^T$  where  $R_1 = (X_{11} + X_{12}) \exp[-D(X_{11} + X_{12})]$  and  $R_2 = (X_{21} + X_{12}) \exp[-D(X_{21} + X_{12})]$ , which consists of correctly specified models. We began by adjusting for the known MGPS.

Next, an analysis was performed in which the GPS was estimated using a generalised linear model for Gamma-distributed data. A univariate or cross-sectional GPS (UGPS) analysis might fail to include information from the previous interval and hence the estimated univariate GPS used would be  $\hat{R}^U = (\hat{R}_1^*, \hat{R}_2^*)^T$  where  $\hat{R}_1^* = (\hat{\theta}_0 + \hat{\theta}_1 X_{11} + \hat{\theta}_2 X_{12}) \exp[-D(\hat{\theta}_0 + \hat{\theta}_1 X_{11} + \hat{\theta}_2 X_{12})]$  and  $\hat{R}_2^* = (\hat{\theta}_0 + \hat{\theta}_1 X_{21}) \exp[-D(\hat{\theta}_0 + \hat{\theta}_1 X_{21})]$ . We compared this with an analysis that adjusts for the estimated multivariate GPS: in this analysis, the response was regressed on treatment and the estimated MGPS

$$\hat{R}^M = (\hat{R}_1, \hat{R}_2)^T \text{ where } \hat{R}_1 = (\hat{\theta}_0 + \hat{\theta}_1 X_{11} + \hat{\theta}_2 X_{12}) \exp[-D(\hat{\theta}_0 + \hat{\theta}_1 X_{11} + \hat{\theta}_2 X_{12})] \text{ and}$$

$$\hat{R}_2 = (\hat{\theta}_0 + \hat{\theta}_1 X_{11} + \hat{\theta}_2 X_{12}) \exp[-D(\hat{\theta}_0 + \hat{\theta}_1 X_{11} + \hat{\theta}_2 X_{12})].$$

**Analyses and results.** We generated 1000 datasets of size 250, 500, 100 and 10 000. Using the true MGPS yielded an APO that lay exactly on the analytically derived dose–response curve (results not shown). The mean and median APO/dose–response curves using the estimated MGPS were also correct, while the UGPS analysis was clearly biased (see Figure 1(a) for results with  $n = 250$ ). The general shape of the UGPS APO was correct; however, the curve fell outside of the confidence bands of the MGPS over part of the range of doses.

### 3.2 Simulation II: Misspecified models

We next considered a more realistic situation, in which the true distribution of the treatment is not known, and so the model is misspecified.

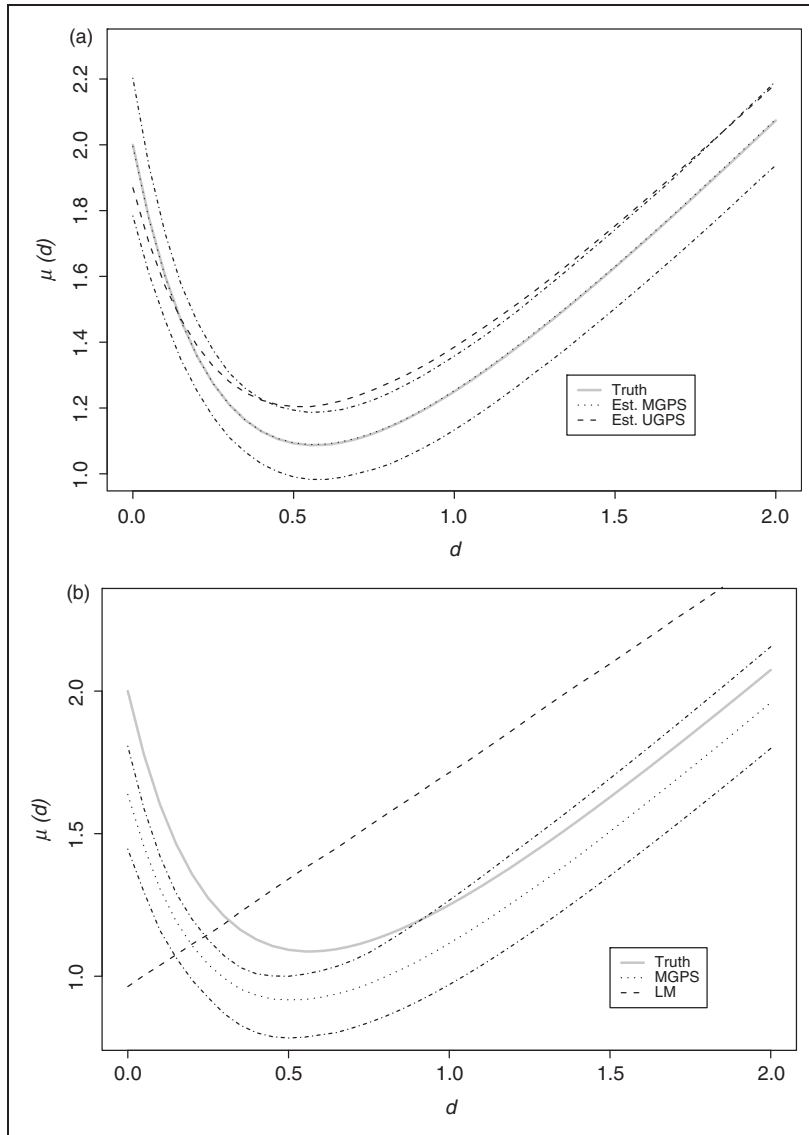
**Data generation.** As before,  $X_{11}$ ,  $X_{12}$  and  $X_{21}$  are each unit exponential and treatment is distributed according to  $D_1 \sim \exp(X_{11} + X_{12})$ ,  $D_2 \sim \exp(X_{21} + X_{12})$ . Now, however, we introduce additional non-linear dependence of the response on the covariates:

$$E[Y_1(d)|X_{11}, X_{12}] = d + (X_{11} + X_{12}) \exp[-d(X_{11} + X_{12})] - 0.3\sqrt{X_{11}X_{12}},$$

$$E[Y_2(d)|X_{21}, X_{12}] = d + (X_{21} + X_{12}) \exp[-d(X_{21} + X_{12})] - 0.3\sqrt{X_{21}X_{12}},$$

where the responses are normal random variables with unit variance. The APO (marginal mean) of the response in each interval is as in the previous simulation.

**Analyses and results.** Two approaches were used to estimate the dose–response relationship: an MGPS analysis and a linear model analysis. In the MGPS analysis, the MGPS was estimated using generalised linear model for Gamma-distributed data. The response was regressed on treatment and the estimated MGPS  $\hat{R} = (\hat{R}_1, \hat{R}_2)^T$  where



**Figure 1.** Simulation results. (a) Simulation Study I: Pointwise median APO adjusted for the estimated MGPS (with pointwise 95% credible interval) and dose-response curve adjusted for the 'univariate' GPS (UGPS). (b) Simulation Study II: Pointwise median dose-response curves, estimated by MGPS-adjusted APO (with pointwise 95% credible interval) and a linear model (LM).

$$\hat{R}_j = (\hat{\theta}_0 + \hat{\theta}_1 X_{j1} + \hat{\theta}_2 X_{12} + \hat{\theta}_3 X_{j1} X_{12}) \times \exp[-D(\hat{\theta}_0 + \hat{\theta}_1 X_{j1} + \hat{\theta}_2 X_{12} + \hat{\theta}_3 X_{j1} X_{12})],$$

for  $j=1,2$  and then the APO was calculated as described in previous sections. In the linear model analysis, the response was modelled as a quadratic function of dose, with adjustment for the confounding variables by the inclusion of linear terms and a first-order interaction between the confounding variables in the regression model. Note that both the MGPS and the linear model

analyses misspecify the relationship between the confounding variables and the response, failing to include the square root of the product of the two confounding variables in the conditional model for the response.

One thousand datasets of size 250, 500, 100 and 10 000 were generated (see Figure 1(b) for results with  $n=250$ ). The mean and median MGPS-adjusted APO curve provides a reasonable approximation of the shape of the true dose–response relationship over most of the range of doses considered. In contrast, the linear model is unable to detect the curvature in the dose–response relationship. Linear model parameter estimates (95% confidence intervals) for  $D$  and  $D^2$  were 0.720 (0.405, 1.034) and 0.016 (−0.020, 0.051), respectively, for  $n=250$ . The dose–response relationship obtained from the linear model analysis and plotted in Figure 1(b) has averaged the estimated dose–response curves over the distribution of the confounding variables, assuming no dependence between covariates and dose.

## 4 Motivating example: MOTAS amblyopia study

We now turn to our motivation for the methodological development of the MGPS, the MOTAS. Amblyopia is the most common childhood vision disorder, and is characterised by reduced visual function in one eye. A standard treatment for the condition is occlusion therapy, that is, patching of the functioning fellow eye. The apparent beneficial effect of occlusion therapy has never been well quantified, partly due to difficulty in the accurate measurement of the occlusion dose. MOTAS<sup>21</sup> was the first clinical study aimed at quantifying the dose–response relationship of occlusion, facilitated by the use of an electronic occlusion dose monitor, consisting of an eye patch with two electrodes attached to its undersurface connected to a data-logger powered by battery from which patch use was read by clinicians at follow-up visits.

The MOTAS design and a full description of the study base have been published previously.<sup>21,22</sup> At study entry, all children who required spectacles entered the *refractive adaptation* phase; the remainder entered the *occlusion* phase directly. Children still considered amblyopic after refractive adaption began occlusion and were prescribed 6 h of occlusion daily. Visual acuity was measured on the logarithm of minimum angle of resolution (logMAR) scale; improvement is indicated by a decrease in logMAR. Visual function and monitored occlusion dose were recorded at approximately 2-week intervals until visual acuity ceased to improve, at which point children exited the study and returned to usual care. A total of 116 children were enrolled in MOTAS; we analyse data of the 68 who took part in the occlusion phase (whether they participated in the refractive adaption phase of the study or not) who, although prescribed 6 h of occlusion daily, received varying occlusion doses because of incomplete concordance. Our notation is as follows: for child  $i$ , the response,  $Y_{ij}$ , is the change in visual acuity during interval  $j$ , and  $D_{ij}$  is the random occlusion dose (in hours) received in interval  $j$ . Intervals are approximately 2 weeks in length, thus a child who concorded perfectly with prescribed treatment would have a dose of 84 h in an interval (i.e. 6 h daily for 14 days). However, children typically did not follow the prescribed occlusion dose, and both higher and lower than prescribed doses were observed.

### 4.1 Applying the MGPS to the MOTAS data

In the study, dose is a continuous variable, but 60 out of 404 (about 15%) of intervals in the occlusion phase had a zero dose. The MGPS model  $f_{D|X}(d|\tilde{x}_{ij}, \beta)$  must acknowledge the mixture nature of the dose distribution, so we assume that, given  $\tilde{X}_{ij} = \tilde{x}_{ij}$ ,

$$D_{ij} \stackrel{\mathcal{L}}{=} \pi(\check{x}_{ij}, \gamma) \mathbb{I}[d = 0] + (1 - \pi(\check{x}_{ij}, \gamma)) \mathbb{I}[d \neq 0] D_{ij}^+, \quad (4)$$

where  $\mathbb{I}[B]$  is the indicator of event  $B$ ,  $D_{ij}^+$  is a strictly positive random variable whose distribution depends on  $\check{X}_{ij} = \check{x}_{ij}$  and  $\beta$ , and  $0 < \pi(\check{x}_{ij}, \gamma) < 1$  is a mixing weight. Estimation in this model is straightforward when a parametric distribution is used for  $D_{ij}^+$ , and any such regression model that induces a balancing property can be used. To estimate  $\gamma$ , we fit a logistic regression model to the binary  $(D_{ij}=0/D_{ij}>0)$  dose data.

Recent work has shown that the binary treatment propensity score should include all confounding variables (i.e. variables predictive of both treatment and outcome) as well as variables that predict outcome, while variables that are purely predictors of treatment should not be included in the model.<sup>23,24</sup> The following covariates therefore included in the MGPS: previous dose, visual acuity at start of interval, age, sex, interval number, length of interval (in days) and amblyopic type (anisometropic, strabismic, mixed). These covariates were used to predict both the probability of having any occlusion at all ( $D_{ij}=0/D_{ij}>0$ ) in a logistic model and the probability of receiving a particular dose (greater than zero) of occlusion in a Weibull model. The MGPS used was

$$\hat{r}(d, \check{x}_{ij}) = \hat{\pi}(\check{x}_{ij}, \hat{\gamma}) \mathbb{I}[d = 0] + (1 - \hat{\pi}(\check{x}_{ij}, \hat{\gamma})) \mathbb{I}[d \neq 0] f(d | \check{x}_{ij}, \hat{\phi}, \hat{\beta}),$$

where  $f(d | \check{x}_{ij}, \phi, \beta)$  is a Weibull density with shape  $\phi$  and scale  $\exp\{\check{x}_{ij}^T \beta\}$ . For the GPS to act as a balancing score, the distribution of  $D_{ij}$  should not depend on  $\check{X}_{ij}$  within strata of  $\hat{r}$ . A graphical check of whether the balancing property was achieved was performed.

As response in the MOTAS is the vector of changes in visual acuity, there is little observed serial correlation in the data. When using a mixture distribution such as (4) for the GPS, it may be that the  $\hat{r}$  values for one component differ substantially from those of the other, so that there are no data in a portion of the space  $\mathcal{D} \times \hat{\mathcal{R}}$  where  $\hat{\mathcal{R}}$  denotes the range of estimated GPS. We account for this explicitly in the model; rather than fitting a model that assumes that the relationship between response and dose and the GPS is the same function in regions of the plane where  $(d, r)$  pairs were observed and in regions where no data was observed, we restrict estimation to a subspace of  $\mathcal{D} \times \hat{\mathcal{R}}$  where data are observed. The observable model for change in visual acuity,  $Y_{ij}(d)$ , is modelled via the expectation

$$E_{Y_{ij}(d)|R_{ij}}[Y_{ij}(d)|R_{ij} = r, \alpha] = \alpha_0 + \mathbb{I}[r < 0.1](\alpha_1 + \alpha_2 d + \alpha_3 r + \alpha_4 d.r), \quad (5)$$

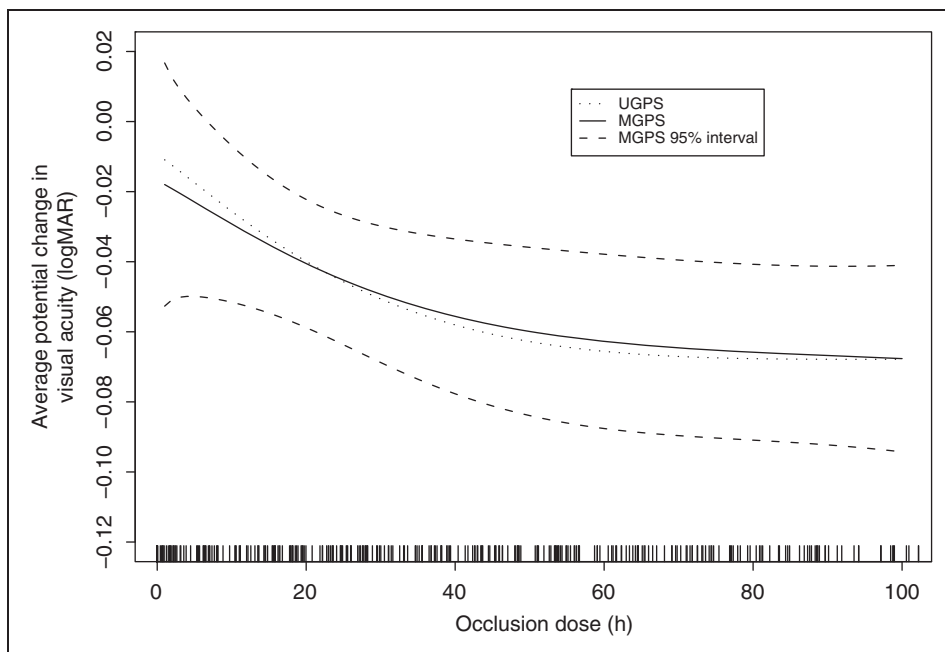
following the model used by Hirano and Imbens.<sup>3</sup> This model can be readily extended to a more flexible or piecewise constant partition model. Here, the inclusion of higher order terms led to only minimal changes in the inferences made. The parameters in the model can be estimated using ordinary least squares or a model that accounts for any remaining within-person correlation, such as a mixed effects model. For the MOTAS data, there was little residual correlation in the longitudinal data for an individual.

## 4.2 Results and comparison with other approaches

The following covariates were included in the MGPS model: occlusion dose in the previous interval, visual acuity at start of interval, age, sex, interval number, length of interval (in days) and amblyopic type. Coefficient estimates (SE) are displayed in Table 1. The distribution of confounding variables

**Table 1.** MOTAS data: Estimates and standard errors for the parameters from the repeated measures (first-order Markov) GPS model: the model comprised a logistic regression for  $D = 0$  versus  $D > 0$  and a Weibull model for positive dose

Model Term	Any dose Est. (SE)	Continuous dose Est. (SE)
Intercept	-2.736 (0.918)	3.490 (0.584)
Previous dose	-0.014 (0.002)	0.003 (0.002)
Visual acuity	0.223 (0.534)	0.562 (0.292)
Age	0.012 (0.009)	-0.001 (0.006)
Sex	-0.392 (0.298)	0.136 (0.165)
Interval number	0.237 (0.053)	0.026 (0.020)
Interval length	0.012 (0.007)	0.003 (0.003)
Type: mixed	-1.901 (0.523)	-0.309 (0.265)
Type: strabismic	-0.239 (0.524)	-0.138 (0.308)



**Figure 2.** MOTAS data: The estimated average potential change in visual acuity (APO) for doses in the range of 1–100 h per interval with pointwise 95% credible interval. The repeated measures GPS (MGPS) APO of section 2.5 is plotted, with a ‘univariate’ GPS APO included for comparison. Observed dose values indicated along the horizontal axis.

such as visual acuity at the start of an interval appeared not to depend on the range of dose within quintiles of the propensity score, indicating that the required balancing property was achieved (results not shown). The observable model in Equation (5) was adopted. Parameters were estimated using a linear mixed effects model with a random intercept and an autoregressive structure to account for any correlation that may exist in the response. Using the model in

Equation (5), we obtain estimates (SE)  $\hat{\alpha}_0 = -0.018(0.008)$ ,  $\hat{\alpha}_1 = -0.002(0.035)$ ,  $\hat{\alpha}_2 = -3.25e - 4(2.93e - 4)$ ,  $\hat{\alpha}_3 = 0.120(3.602)$  and  $\hat{\alpha}_4 = -0.070(0.083)$ , respectively; the standard deviation of the random intercept term was  $2.17e - 06$  and the autoregressive correlation was 0.077.

A plot of the dose-response curve is presented in Figure 2, with numerical values of the estimated APO presented in Table 2. The plot indicates that the direct effect of dose on visual acuity, when confounding between dose and the responses is adjusted for using the GPS approach, is appreciable; the average potential effect on change in visual acuity measurement  $Y_{ij}$  is significantly negative (corresponding to vision improvement) over the entire range of positive doses considered.

Using the model in Equation (5) with a 'univariate' GPS strategy, i.e. one in which previous dose was not included in the dose density function, we obtain least-squares estimates (SE)  $\hat{\alpha}_0 = -0.022(0.007)$ ,  $\hat{\alpha}_1 = 0.010(0.033)$ ,  $\hat{\alpha}_2 = -2.86 - 04(2.94e - 04)$ ,  $\hat{\alpha}_3 = 0.219(3.382)$  and  $\hat{\alpha}_4 = -0.107(0.078)$ , respectively. The APO obtained is very similar to that obtained using the MGPS, indicating that previous dose did little to confound the association between most recent dose and change in visual acuity.

Finally, the average dose effect estimated by a mixed effects model analysis yields an intercept further from zero than the GPS analyses, and does not find any evidence of curvature in the dose-response relationship. The GPS results are in closer agreement with the ophthalmological belief that visual acuity will not improve spontaneously in the absence of occlusion. Also, the GPS APOs suggest a plateau, or a saturation, of the effect of occlusion in an interval at about 80 h. This indicates that children may not exhibit a clinically meaningful improvement in visual acuity with more than, on average, 6 h of occlusion per day over a 2-week period. This is biologically plausible, as physical changes to the amblyopic eye that can occur in a fixed time period are likely limited by physiological processes.

## 5 Discussion

In the use of propensity scores for binary treatments, there has been considerable attention devoted of late to the use of propensity scores to design observational studies,<sup>5</sup> or to estimate average effects only when the experimental treatment assignment assumption is satisfied.<sup>25</sup> That is, an average causal effect can only be identified without further assumptions (and hence should only be estimated using the entire sample of observed data) if all individuals in the population have a positive probability of receiving both the active and control treatments. A straightforward check of this can be performed by examining the distribution of propensity scores among the treated and

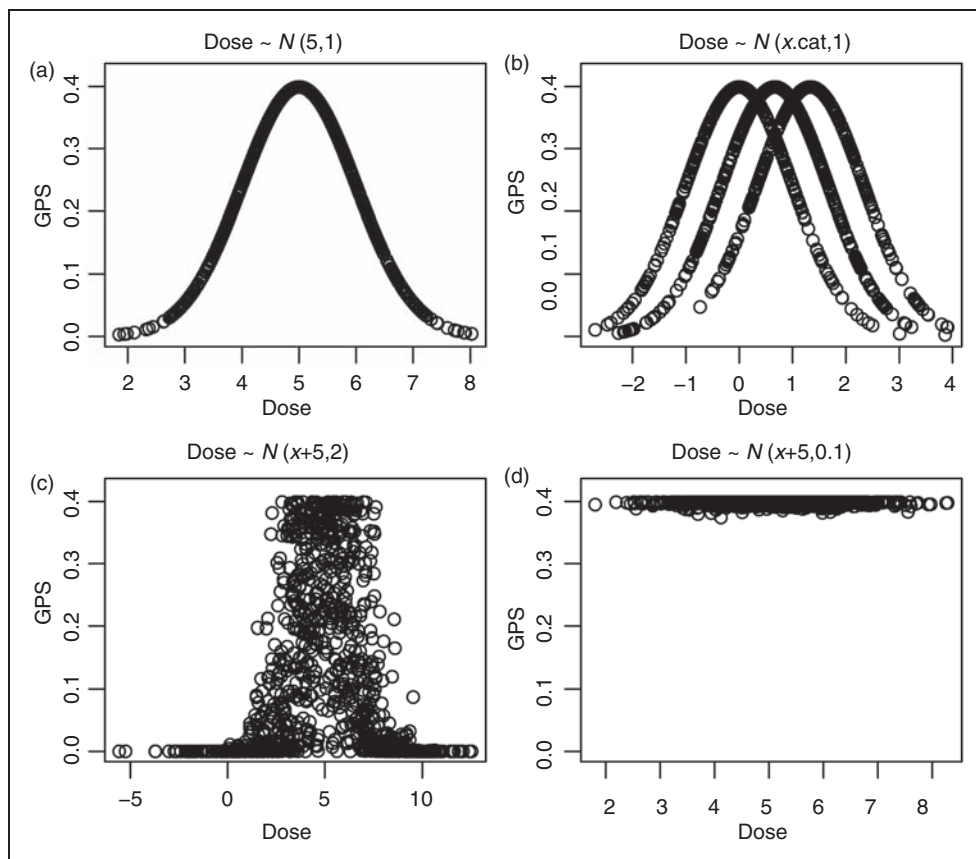
**Table 2.** MOTAS data: Summaries of the APO (on the logMAR scale) from a multivariate GPS model analysis for changing dose amount per interval: 5000 bootstrap samples

Quantile	Dose (h)				
	10	25	50	75	100
2.5	-0.044	-0.055	-0.085	-0.102	-0.114
25	-0.028	-0.043	-0.070	-0.086	-0.096
50	-0.021	-0.036	-0.061	-0.077	-0.086
75	-0.013	-0.030	-0.053	-0.068	-0.077
97.5	0.002	-0.018	-0.039	-0.052	-0.060



the untreated to see whether these are completely overlapping;<sup>26</sup> however, more sophisticated approaches have also been proposed.<sup>27,28</sup>

When treatment doses are continuous, there is no simple analog to plotting the GPS values against dose levels to look for overlap. Consider, for illustration, a setting in which there is a single covariate  $X$  that may confound the relationship between dose and the outcome of interest. In Figure 3, the relationship between a normally distributed dose and the GPS is plotted under four scenarios of dependence of dose on the covariate. In the first scenario, dose level does not depend on  $X$ , while in the second, the mean dose depends on categorised values of  $X$ . The final two scenarios allow the mean dose level to depend on  $X$  in a linear fashion with different degrees of strength of the relationship (correlations of approximately 0.34 and 0.99). It is clear that, for normally-distributed doses, it is possible to distinguish to some extent between the different degrees of dependence of dose on a single covariate. Furthermore, by examining the situation where dose depends on a categorical variable, an approximate guideline presents itself for deciding the dose levels within which it is



**Figure 3.** Scatter plots of the true GPS by dose when dose level is normally distributed. (a) Dose level does not depend on  $X$ , (b) dose depends on categorised values of  $X$ , (c) and (d) dose level depends linearly on  $X$  (smaller variance in the right panel).

reasonable to estimate the APO. In particular, we wish to find a range of doses that is no wider than four standard deviations (so as to contain approximately 95% of the distribution) and that contains the peaks of the dose density curves.

It may be the case that the peaks of the density curves are not sufficiently close together that all can be contained in an interval that is four standard deviations in length. In such a case – for example, if the dependence of dose on covariates is very high – it may be necessary to estimate APOs in subpopulations that differ by dosage level, rather than attempt to construct a population average dose–response estimate over the entire range of doses  $\mathcal{D}$  which would rely on smoothing over regions of the covariate  $\times$  dose space in which there are little or no data.

Unfortunately, it may be more difficult to identify whether there is sufficient variability in doses over the covariate range, or over which range to estimate the dose–response curve, simply by plotting the GPS scores against the dose values. For example, if the dose distribution is asymmetric, the doses in which to estimate the APO should be selected to maximise the number of observations contained in a four standard deviation range. This will help to reduce the loss of power that results from the restriction of the range and the consequent elimination of data points lying outside of that range. Research into whether all dose levels could have been received by all individuals as described by their covariate types is required.

Another important point of consideration in any propensity score analysis is the issue of model choice. A number of authors<sup>23,24,29</sup> have considered the case of binary treatment propensity scores, and concluded that including in the propensity score model all variables that are causes of both the response and the treatment (confounding variables), and all variables that are predictors of the response only, improves the performance of the estimators, while the inclusion of variables that are causes of treatment only is not helpful and leads to an increase in variance. As noted here, GPS adjustment simplifies the model specification for the response, which can be modelled only as a function of dose and a one-dimensional summary of the covariates, namely the GPS. However, it may be the case that the dose–response curve is modified by some covariate(s). In such a case – for example, where treatment interacts with sex – it may be preferable to include sex and the dose-by-sex interaction directly in the response model rather than including sex in the dose model, particularly if dose–response profiles for each sex separately were desired (rather than the marginal population-averaged dose–response profile).

We have extended the GPS methodology to the repeated measures setting to cope with situations where treatment doses received in different intervals are correlated and response may depend on doses in current and earlier intervals. In a longitudinal study of dose–response, full compliance is the exception rather than the expected. To estimate the dose–response relationship with confidence, modelling potentially confounding relationships flexibly is key. The GPS is under-used, yet the approach provides a tractable and flexible option of analysis, and can be adapted to analyse any number of complex dosing strategies – including multi-interval treatments.

## Acknowledgements

Both authors acknowledge funding through the Natural Sciences and Engineering Research Council of Canada (NSERC). MOTAS was supported by The Guide Dogs for the Blind Association, UK.

## References

1. Arjas E and Parner J. Causal reasoning from longitudinal data. *Scandinavian Journal of Statistics* 2004; **31**: 171–87.
2. Imbens GW. The role of the propensity score in estimating dose-response functions. *Biometrika* 2000; **87**: 706–10.

3. Hirano K and Imbens GW. The propensity score with continuous treatments. In: Gelman A and Meng XL (eds) *Applied bayesian modeling and causal inference from incomplete-data perspectives*. Oxford, UK: Wiley, 2004, pp.73–84.
4. Flores C. *Estimation of dose-response functions and optimal doses with a continuous treatment*. Technical report. University of Miami, 2004.
5. Rubin DB. The design versus the analysis of observational studies for causal effects: Parallels with the design of randomised trials. *Statistics in Medicine* 2007; **26**: 20–36.
6. Follman DA. On the effect of treatment among would-be treatment compliers: an analysis of the multiple risk factor intervention trial. *Journal of the American Statistical Association* 2000; **95**: 1101–9.
7. Frangakis CE and Rubin DB. Addressing complications of intention-to-treat analysis in the combined presence of all-or-none treatment-compliance and subsequent missing outcomes. *Biometrika* 1999; **86**: 365–79.
8. Frangakis CE and Rubin DB. Principal stratification in causal inference. *Biometrics* 2002; **58**: 21–9.
9. Joffe MM, Ten Have TR and Brensinger C. The compliance score as a regressor in randomised trials. *Biostatistics* 2003; **4**: 327–40.
10. Neugebauer R and van der Laan M. G-computation estimation of nonparametric causal effects on time-dependent mean outcomes in longitudinal studies. *Computational Statistics & Data Analysis* 2006; **51**: 1676–97.
11. Robins JM, Hernán MA and Brumback B. Marginal structural models and causal inference in epidemiology. *Epidemiology* 2000; **11**: 550–60.
12. Hernán MA, Brumback B and Robins JM. Marginal structural models to estimate the causal effect of zidovudine on the survival of HIV-positive men. *Epidemiology* 2000; **11**: 561–70.
13. Hernán MA, Lanoy E, Costagliola D and Robins JM. Comparison of dynamic treatment regimes via inverse probability weighting. *Basic & Clinical Pharmacology & Toxicology* 2006; **98**: 237–42.
14. Bembom O and van der Laan MJ. Analyzing sequentially randomised trials based on causal effect models for realistic individualised treatment rules. *COBRA (Collection of Biostatistics Research Archive)*; 2007.
15. van der Laan MJ and Petersen ML. Causal effect models for realistic individualized treatment and intention to treat rules. *International Journal of Biostatistics* 2007; **3**(1), Article 6.
16. Robins J, Orellana L and Rotnitzky A. Estimation and extrapolation of optimal treatment and testing strategies. *Statistics in Medicine* 2008; **27**(23): 4678–721.
17. Imai K and Van Dyk DA. Causal inference with general treatment regimes: Generalizing the propensity score. *Journal of the American Statistical Association* 2004; **99**(467): 854–66.
18. Robins JM. Causal inference from complex longitudinal data. In: Berkane M (ed.) *Latent variable modeling and applications to causality*. New York: Springer-Verlag, 1997, pp.69–117.
19. Rubin DB. Bayesian inference for causal effects: The role of randomization. *Annals of Statistics* 1978; **6**: 34–58.
20. Rosenbaum PR and Rubin DB. The central role of the propensity score in observational studies for causal effects. *Biometrika* 1983; **70**: 41–55.
21. Stewart CE, Moseley MJ, Stephens DA and Fielder AR. Treatment dose–response in amblyopia therapy: The Monitored Occlusion Treatment of Amblyopia Study (MOTAS). *Investigations in Ophthalmology and Visual Science* 2004; **45**: 3048–54.
22. Stewart CE, Fielder AR, Stephens DA and Moseley MJ. Design of the Monitored Occlusion Treatment of Amblyopia Study (MOTAS). *British Journal of Ophthalmology* 2002; **86**: 915–9.
23. Brookhart MA, Schneeweiss S, Rothman KJ, Glynn RJ, Avorn J and Sturmer T. Variable selection for propensity score models. *American Journal of Epidemiology* 2006; **163**: 1149–56.
24. Austin PC, Grootendorst P and Anderson GM. A comparison of the ability of different propensity score models to balance measured variables between treated and untreated subjects: A Monte Carlo study. *Statistics in Medicine* 2007; **26**: 734–53.
25. Bembom O and van der Laan M. A practical illustration of the importance of realistic individualised treatment rules in causal inference. *Electronic Journal of Statistics* 2007; **1**: 574–96.
26. Glynn RJ, Schneeweiss S and Stürmer T. Indications for propensity scores and review of their use in pharmacoepidemiology. *Basic & Clinical Pharmacology & Toxicology* 2006; **98**: 253–9.
27. Wang Y, Petersen ML, Bangsberg D and van der Laan MJ. Diagnosing bias in the inverse probability of treatment weighted estimator resulting from violation of experimental treatment assignment. *COBRA (Collection of Biostatistics Research Archive)*; 2006.
28. Imai K, King G and Stuart EA. Misunderstandings between experimentalists and observationalists about causal inference. *Journal of the Royal Statistical Society, Series A* 2008; **171**: 481–502.
29. Judkins DR, Morganstein D, Zador P, Piesse A, Barrett B and Mukhopadhyay P. Variable selection and raking in propensity scoring. *Statistics in Medicine* 2007; **26**: 1022–33.

## Appendix: Theoretical properties of the MGPS

**Proof of Theorem 1, Weak unconfoundedness given the MGPS:** Sequential weak unconfoundedness implies that for all  $d \in \mathcal{D}_j$ ,  $Y_{ij}(d) \perp D_{ij} | \check{X}_{ij}$ , that is, for each dose  $d$ ,  $Y_{ij}(d)$  and  $D_{ij}$  are conditionally independent given  $\check{X}_{ij}$ . To establish the result we need to show that

$$f_{D_{ij}|R_{ij}, Y_{ij}(d)}(d | r(d, \check{X}_{ij}), y_{ij}(d)) = f_{D_{ij}|R_{ij}}(d | r(d, \check{X}_{ij})).$$

Consider the random quantity,  $R_{ij} = r(d, \check{X}_{ij})$  where  $r(d, \check{X}_{ij}) = f_{D_{ij}|X_{ij}}(d | \check{X}_{ij})$ , defined for fixed  $d$ . Denoting by  $f$  the density function for the relevant random variables, we have

$$\begin{aligned}
f_{D_{ij}|R_{ij}}(d|r(d, \check{x}_{ij})) &= \int_{\mathcal{X}_{d,j}} f_{D_{ij}, \check{X}_{ij}|R_{ij}}(d, \check{x}|r(d, \check{x}_{ij})) d\check{x} \\
&= \int_{\mathcal{X}_{d,j}} f_{D_{ij}|\check{X}_{ij}, R_{ij}}(d|\check{x}, r(d, \check{x}_{ij})) f_{\check{X}_{ij}|R_{ij}}(\check{x}|r(d, \check{x}_{ij})) d\check{x}
\end{aligned}$$

where  $\mathcal{X}_{d,j} \subset \mathcal{X}_j$  is the set of solutions  $\check{x}$  of the equation  $r(d, \check{x}) = r(d, \check{x}_{ij})$ . Now, in  $\mathcal{X}_{d,j}$

$$f_{D_{ij}|\check{X}_{ij}, R_{ij}}(d|\check{x}, r(d, \check{x}_{ij})) \equiv f_{D_{ij}|\check{X}_{ij}}(d|\check{x}_{ij}) = r(d, \check{x}_{ij}),$$

as for fixed  $d$  and  $\check{x}$  in  $\mathcal{X}_{d,j}$ ,  $r(d, \check{x}) = r(d, \check{x}_{ij})$  is completely defined. Thus

$$\begin{aligned}
f_{D_{ij}|R_{ij}}(d|r(d, \check{x}_{ij})) &= \int_{\mathcal{X}_{d,j}} r(d, \check{x}_{ij}) f_{\check{X}_{ij}|R_{ij}}(\check{x}|r(d, \check{x}_{ij})) d\check{x} \\
&= r(d, \check{x}_{ij}) \int_{\mathcal{X}_{d,j}} f_{\check{X}_{ij}|R_{ij}}(\check{x}|r(d, \check{x}_{ij})) d\check{x} = r(d, \check{x}_{ij}) = f_{D_{ij}|\check{X}_{ij}}(d|\check{x}_{ij})
\end{aligned}$$

for an archetypal  $\check{x}_{ij}$  in  $\mathcal{X}_{d,j}$ . Similarly, by weak unconfoundedness,

$$\begin{aligned}
f_{D_{ij}|R_{ij}, Y_{ij}(d)}(d|r(d, \check{x}_{ij}), y_{ij}(d)) &= \int_{\mathcal{X}_{d,j}} f_{D_{ij}|\check{X}_{ij}, R_{ij}, Y_{ij}(d)}(d|\check{x}, r(d, \check{x}_{ij}), y_{ij}(d)) \\
&\quad f_{\check{X}_{ij}|R_{ij}, Y_{ij}(d)}(\check{x}|r(d, \check{x}_{ij}), y_{ij}(d)) d\check{x} \\
&= \int_{\mathcal{X}_{d,j}} f_{D_{ij}|\check{X}_{ij}}(d|\check{x}) f_{\check{X}_{ij}|R_{ij}, Y_{ij}(d)}(\check{x}|r(d, \check{x}_{ij}), y_{ij}(d)) d\check{x} \\
&= \int_{\mathcal{X}_{d,j}} r(d, \check{x}_{ij}) f_{\check{X}_{ij}|R_{ij}, Y_{ij}(d)}(\check{x}|r(d, \check{x}_{ij}), y_{ij}(d)) d\check{x} \\
&= r(d, \check{x}_{ij}) = f_{D_{ij}|\check{X}_{ij}}(d|\check{x}_{ij})
\end{aligned}$$

for an archetypal  $\check{x}_{ij}$  in  $\mathcal{X}_{d,j}$ . Thus, for all  $d$ , we have weak unconfoundedness given  $R_{ij} = r(d, \check{x}_{ij})$ .

**Proof of Theorem 2, Bias removal of the MGPS procedure:** Consider the conditional distribution of potential response  $Y_{ij}(d)$  given  $D_{ij} = d$  and  $R_{ij} = r_{ij} = r(d, \check{x}_{ij})$ , for fixed  $d \in \mathcal{D}$ . By conditional probability and Theorem 1 above, we have

$$\begin{aligned}
f_{Y_{ij}(d)|D_{ij}, R_{ij}}(y_{ij}(d)|d, r(d, \check{x}_{ij})) &= \frac{f_{Y_{ij}(d)|R_{ij}}(y_{ij}(d)|r(d, \check{x}_{ij})) f_{D_{ij}|Y_{ij}(d), R_{ij}}(d|y_{ij}(d), r(d, \check{x}_{ij}))}{f_{D_{ij}|R_{ij}}(d|r(d, \check{x}_{ij}))} \\
&= \frac{f_{Y_{ij}(d)|R_{ij}}(y_{ij}(d)|r(d, \check{x}_{ij})) f_{D_{ij}|R_{ij}}(d|r(d, \check{x}_{ij}))}{f_{D_{ij}|R_{ij}}(d|r(d, \check{x}_{ij}))} \\
&= f_{Y_{ij}(d)|R_{ij}}(y_{ij}(d)|r(d, \check{x}_{ij})).
\end{aligned}$$

However, by definition and Theorem 1,

$$E[Y_{ij}(d)|R_{ij} = r(d, \check{x}_{ij})] = \beta(d, r(d, \check{x}_{ij}))$$

Thus, by iterated expectation, noting that  $E[Y(d)] \equiv E[Y_{ij}(d)]$ ,

$$\mu(d) \equiv E[Y_{ij}(d)] = E_{R_{ij}}[E[Y_{ij}(d)|R_{ij} = r(d, \check{x}_{ij})]] \equiv E_{\check{X}_{ij}}[\beta(d, r(d, \check{x}_{ij}))].$$