

LONDON
SCHOOL of
HYGIENE
& TROPICAL
MEDICINE



Diazordaz, K; Franchini, AJ; Grieve, R (2017) Methods for estimating complier average causal effects for cost-effectiveness analysis. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*. n/a-n/a. ISSN 1467-985X DOI: <https://doi.org/10.1111/rssa.12294>

Downloaded from: <http://researchonline.lshtm.ac.uk/4645614/>

DOI: [10.1111/rssa.12294](https://doi.org/10.1111/rssa.12294)

Usage Guidelines

Please refer to usage guidelines at <http://researchonline.lshtm.ac.uk/policies.html> or alternatively contact researchonline@lshtm.ac.uk.

Available under license: <http://creativecommons.org/licenses/by-nc-nd/2.5/>

Methods for estimating complier-average causal effects for cost-effectiveness analysis

K. DiazOrdaz, A. J. Franchini, R. Grieve

London School of Hygiene and Tropical Medicine

†E-mail: karla.diaz-ordaz@lshtm.ac.uk

Summary. In Randomised Controlled Trials (RCT) with treatment non-compliance, instrumental variable approaches are used to estimate complier average causal effects. We extend these approaches to cost-effectiveness analyses, where methods need to recognise the correlation between cost and health outcomes. We propose a Bayesian full likelihood (BFL) approach, which jointly models the effects of random assignment on treatment received and the outcomes, and a three-stage least squares (3sls) method, which acknowledges the correlation between the endpoints, and the endogeneity of the treatment received. This investigation is motivated by the REFLUX study, which exemplifies the setting where compliance differs between the RCT and routine practice. A simulation is used to compare the methods performance. We find that failure to model the correlation between the outcomes and treatment received correctly can result in poor CI coverage and biased estimates. By contrast, BFL and 3sls methods provide unbiased estimates with good coverage.

Keywords: non-compliance, instrumental variables, bivariate outcomes, cost-effectiveness

1. Introduction

Non-compliance is a common problem in Randomised Controlled Trials (RCTs), as some participants depart from their randomised treatment, by for example switching from the experimental to the control regimen. An unbiased estimate of the effectiveness of treatment assignment can be obtained by reporting the intention-to-treat (ITT) estimand. In the presence of non-compliance, a complimentary estimand of interest is the causal effect of treatment received. Instrumental variable (IV) methods can be used to obtain the complier average causal effect (CACE), as long as random assignment meets the IV criteria for identification (Angrist et al., 1996). An established approach to IV estimation is two-stage least squares (2sls), which provides consistent estimates of the CACE when the outcome measure is continuous, and non-compliance is binary (Baiocchi et al., 2014).

†*Address for correspondence:* Department of Medical Statistics, London School of Hygiene and Tropical Medicine, Keppel Street, London, WC1E 7HT, UK.

Cost-effectiveness analyses (CEA) are an important source of evidence for informing clinical decision-making and health policy. CEA commonly report an ITT estimand, *i.e.* the relative cost-effectiveness of the intention to receive the intervention (NICE, 2013). However, policy-makers may require additional estimands, such as the relative cost-effectiveness for compliers. For example, CEAs of new therapies for end-stage cancer, are required to estimate the cost-effectiveness of treatment receipt, recognising that patients may switch from their randomised allocation following disease progression. Alternative estimates such as the CACE may also be useful when levels of compliance in the RCT differ to those in the target population, or where intervention receipt, rather than the intention to receive the intervention, is the principal cost driver. Methods for obtaining the CACE for univariate survival outcomes have been exemplified before (Latimer *et al.*, 2014), but approaches for obtaining estimates that adequately adjust for non-adherence in CEA more generally, have received little attention. This has been recently identified as a key area where methodological development is needed (Hughes *et al.*, 2016).

The context of trial-based CEA highlights an important complexity that arises with multivariate outcomes more widely, in that, to provide accurate measures of the uncertainty surrounding a composite measure of interest, for example the incremental net monetary benefit (INB), it is necessary to recognise the correlation between the endpoints, in this case, cost and health outcomes (Willan *et al.*, 2003; Willan, 2006). Indeed, when faced with non-compliance, and the requirement to estimate a causal effect of treatment on cost-effectiveness endpoints, some CEA resort to per protocol (PP) analyses (Brilleman *et al.*, 2015), which exclude participants who deviate from treatment. As non-compliance is likely to be associated with prognostic variables, only some of which are observed, PP analyses are liable to provide biased estimates of the causal effect of the treatment received.

This paper develops novel methods for estimating CACE in CEA that use data from RCTs with non-compliance. First, we propose using the three stage least squares (3sls) method (Zellner and Theil, 1962), which allows the estimation of a system of simultaneous equations with endogenous regressors. Next, we consider a bivariate version of the ‘unadjusted Bayesian’ models previously proposed for Mendelian randomisation (Burgess and Thompson, 2012), which simultaneously estimate the expected treatment received as a function of random allocation, and the mean outcomes as a linear function of the expected treatment received. Finally, we develop a Bayesian full likelihood approach (BFL), whereby the outcome variables and the treatment received are jointly modelled as dependent on the random assignment. This is an extension to the multivariate case of what is known in the econometrics literature as the IV unrestricted reduced form (Kleibergen and Zivot, 2003).

The aim of this paper is to present and compare these alternative approaches. The problem is illustrated in Section 2 with the REFLUX study, a multicentre RCT and CEA that contrasts laparoscopic surgery with medical management for patients with Gastro-Oesophageal Reflux Disease (GORD). Section 3 introduces the assumptions and

methods for estimating CACEs. Section 4 presents a simulation study used to assess the performance of the alternative approaches, which are then applied to the case study in Section 5. We conclude with a Discussion (Section 6), where we consider the findings from this study in the context of related research.

2. Motivating example: Cost-effectiveness analysis of the REFLUX study

The REFLUX study was a UK multicentre RCT with a parallel design, in which patients with moderately severe GORD, were randomly assigned to medical management or laparoscopic surgery (Grant et al., 2008, 2013).

The RCT randomised 357 participants (178 surgical, 179 medical) from 21 UK centres. An observational preference based study was conducted alongside it, which involved 453 preference participants (261 surgical, 192 medical).

For the cost-effectiveness analysis within the trial, individual resource use (costs in £ sterling) and health-related quality of life (HRQoL), measured using EQ5D (3 levels), were recorded annually for up to 5 years. The HRQoL data were used to adjust life years and present quality-adjusted life years (QALYs) over the follow-up period (Grant et al., 2013).[‡] As is typical, the costs were right-skewed. Table 1 reports the main characteristics of the data set.

The original CEA estimated the linear additive treatment effect on mean costs and health outcomes (QALYs). The primary analysis used a system of seemingly unrelated regression equations (SURs) (Zellner, 1962; Willan et al., 2004), adjusting for baseline HRQoL EQ5D summary score (denoted by EQ5D₀). The SURs can be written for cost Y_{1i} and QALYs Y_{2i} , as follows

$$\begin{aligned} Y_{1i} &= \beta_{0,1} + \beta_{1,1}\text{treat}_i + \beta_{1,2}\text{EQ5D}_{0i} + \epsilon_{1i} \\ Y_{2i} &= \beta_{0,2} + \beta_{1,2}\text{treat}_i + \beta_{2,2}\text{EQ5D}_{0i} + \epsilon_{2i} \end{aligned} \quad (1)$$

where $\beta_{1,1}$ and $\beta_{1,2}$ represent the incremental costs and QALYs respectively. The error terms are required to satisfy $E[\epsilon_{1i}] = E[\epsilon_{2i}] = 0$, $E[\epsilon_{ki}\epsilon_{k'i}] = \sigma_{kk'}$, $E[\epsilon_{ki}\epsilon_{k'j}] = 0$, for $k, k' \in \{1, 2\}$, and for $i \neq j$. Rather than assuming that the errors are drawn from a bivariate normal distribution, estimation is usually done by the feasible generalized least squares (FGLS) method[§]. This is a two-step method where, in the first step, we run ordinary least squares estimation for equation (1). In the second step, residuals from the first step are used as estimates of the elements $\sigma_{kk'}$ of the covariance matrix, and this estimated covariance structure is then used to re-estimate the coefficients in equation (1) (Zellner, 1962; Zellner and Huang, 1962).

In addition to reporting incremental costs and QALYs, CEA often report the incremental cost-effectiveness ratio (ICER), which is defined as the ratio of the incremental costs per incremental QALY, and the incremental net benefit (INB), defined as

[‡]There was no administrative censoring.

[§]If we are prepared to assume the errors are bivariate normal, estimation can proceed by maximum likelihood.

$INB(\lambda) = \lambda\beta_{1,2} - \beta_{1,1}$, where λ represents the decision-makers' *willingness to pay* for a one unit gain in health outcome. Thus the new treatment is cost-effective if $INB > 0$. For a given λ , the standard error of INB can be estimated from the estimated increments $\hat{\beta}_{1,1}$ and $\hat{\beta}_{1,2}$, together with their standard errors and their correlation following the usual rules for the variance of a linear combination of two random variables. The willingness to pay λ generally lies in a range, so it is common to compute the estimated value of INB for various values of λ . In REFLUX, the reported INB was calculated using $\lambda = \pounds 30000$, which is within the range of cost-effectiveness thresholds used by the UK National Institute for Health and Care Excellence (NICE, 2013).

The original ITT analysis concluded that, compared to medical management, the arm assigned to surgery had a large gain in average QALYs, at a small additional cost and was relatively cost-effective with a positive mean INB , albeit with 95% confidence intervals (CI) that included zero. However, these ITT results cannot be interpreted as a causal effect of the treatment, since within one year of randomisation, 47 of those randomised to surgery switched and received medical management, while in the medical treatment arm, 10 received surgery. The reported reasons for not having the allocated surgery were that in the opinion of the surgeon or the patient, the symptoms were not “sufficiently severe” or the patient was judged unfit for surgery (e.g. overweight). The preference-based observational study conducted alongside the RCT reported that in routine clinical practice, the corresponding proportion who switched from an intention to have surgery and received medical management, was relatively low (4%), with a further 2% switching from control to intervention (Grant et al., 2013). Since the percentage of patients who switched in the RCT was higher than in the target population and the costs of the receipt of surgery are relatively large, there was interest in reporting a causal estimate of the intervention. Thus, the original study also reported a PP analysis on complete-cases, adjusted for baseline EQ5D₀, which resulted in an ICER of $\pounds 7263$ per additional QALY (Grant et al., 2013). This is not an unbiased estimate of the causal treatment effect, so in Section 5, we re-analyse the REFLUX dataset to obtain a CACE of the cost-effectiveness outcomes, recognising the joint distribution of costs and QALYs, using the methods described in the next section.

3. Complier Average Causal effects with bivariate outcomes

We begin by defining more formally our estimands and assumptions. Let Y_{1i} and Y_{2i} be the continuous bivariate outcomes, and Z_i and D_i the binary random treatment allocation and treatment received respectively, corresponding to the i -th individual. The bivariate endpoints Y_{1i} and Y_{2i} belong to the same individual i , and thus are correlated. We assume that there is an unobserved confounder U , which is associated with the treatment received and either or both of the outcomes. From now on, we will assume that the **(i) Stable Unit Treatment Value Assumption (SUTVA)** holds: the potential outcomes of the i -th individual are unrelated to the treatment status of all other individuals (known as *no*

interference), and that for those who actually received treatment level z , their observed outcome is the potential outcome corresponding to that level of treatment.

Under SUTVA, we can write the potential treatment received by the i -th subject under the random assignment at level $z_i \in \{0, 1\}$ as $D_i(z_i)$. Similarly, $Y_{\ell i}(z_i, d_i)$ with $\ell \in \{1, 2\}$ denotes the corresponding potential outcome for endpoint ℓ , if the i -th subject were allocated to level z_i of the treatment and received level d_i . There are four potential outcomes. Since each subject is randomised to one level of treatment, only one of the potential outcomes per endpoint ℓ , is observed, i.e. $Y_{\ell i} = Y_{\ell i}(z_i, D_i(z_i)) = Y_i(z_i)$.

The CACE for outcome ℓ can now be defined as

$$\theta_{\ell} = E \left[\{Y_{\ell i}(1) - Y_{\ell i}(0)\} \{D_i(1) - D_i(0) = 1\} \right]. \quad (2)$$

In addition to (i) SUTVA, the following assumptions are sufficient for identification of the CACE, (Angrist et al., 1996):

(ii) Ignorability of the treatment assignment: Z_i is independent of unmeasured confounders (conditional on measured covariates) and the potential outcomes $Z_i \perp\!\!\!\perp U_i, D_i(0), D_i(1), Y_i(0), Y_i(1)$.

(iii) The random assignment predicts treatment received: $Pr\{D_i(1) = 1\} \neq Pr\{D_i(0) = 1\}$.

(iv) Exclusion restriction: The effect of Z on Y_{ℓ} must be via an effect of Z on D ; Z cannot affect Y_{ℓ} directly.

(v) Monotonicity: $D_i(1) \geq D_i(0)$.

The CACE can now be identified from equation (2) without any further assumptions about the unobserved confounder; in fact, U can be an effect modifier of the relationship of D and Y (Didilez et al., 2010).

In the REFLUX study, the assumptions concerning the random assignment, (ii) and (iii), are justified by design. The exclusion restriction assumption seems plausible for the costs, since the costs of surgery are only incurred if the patient actually has the procedure. We argue that it is also plausible that it holds for QALYs, as the participants did not seem to have a preference for either treatment, thus making the psychological effects of knowing to which treatment one has been allocated minimal. The monotonicity assumption rules out the presence of defiers. It seems fair to assume that there are no participants who would refuse the REFLUX surgery when randomised to it, but who would receive surgery when randomised to receive medical management. Equation (2) implicitly assumes that receiving the intervention has the same average effect in the linear scale, regardless of the level of Z and U . This average is however across different ‘versions’ of the intervention, as the trial protocol did not prescribe a single surgical procedure, but allowed for the surgeon to choose their preferred laparoscopy method, as would be the case in routine clinical practice.

Since random allocation, Z , satisfies assumptions (ii)-(iv), we say it is an instrument (or instrumental variable) for D . For binary instrument, the simplest method of estimation

of equation (2) in the IV framework is the Wald estimator (Angrist et al., 1996):

$$\hat{\theta}_{\ell,IV} = \frac{E(Y|Z = 1) - E(Y|Z = 0)}{E(D|Z = 1) - E(D|Z = 0)}$$

Typically, estimation of these conditional expectations proceeds via an approach known as two-stage least squares (2sls). The first stage fits a linear regression to treatment received on treatment assigned. Then, in a second stage, a regression model for the outcome on the predicted treatment received is fitted:

$$\begin{aligned} D_i &= \alpha_0 + \alpha_1 Z_i + \omega_{1i} \\ Y_{\ell i} &= \beta_0 + \beta_{IV} \hat{D}_i + \omega_{2i} \end{aligned} \quad (3)$$

where $\hat{\beta}_{IV}$ is an estimator for θ_{ℓ} . Covariates can be used, by including them in both stages of the model. To obtain the correct standard errors for the 2sls estimator, it is necessary to take into account the uncertainty about the first stage estimates. The asymptotic standard error for the 2sls CACE is given in Imbens and Angrist (1994), and implemented in commonly used software packages.

OLS estimation produces first-stage residuals ω_{1i} that are uncorrelated with the instrument, and this is sufficient to guarantee that the 2sls estimator is consistent for the CACE (Angrist et al., 2008). Therefore, we restrict our attention here to models where the first-stage equation is linear, even though the treatment received is binary. ¶

A key issue for settings such as CEA where there is interest in estimating the CACE for bivariate outcomes, is that 2sls as implemented in most software packages can only be readily applied to univariate outcomes. Ignoring the correlation between the two endpoints is a concern for obtaining standard errors of composite measures of the outcomes, e.g. INB, as this requires accurate estimates of the covariance between the outcomes of interest (e.g. costs and QALYs).

A simple way to address this problem would be to apply 2sls directly to the composite measure, i.e. a net-benefit two-stage regression approach (Hoch et al., 2006). However, it is known that net benefit regression is very sensitive to outliers, and to distributional assumptions (Willan et al., 2004), and has been recently shown to perform poorly when these assumptions are thought to be violated (Mantopoulos et al., 2016). Moreover, such net benefit regression is restrictive, in that it does not allow separate covariate adjustment for each of the component outcomes, (e.g. baseline HRQoL, for the QALYs as opposed to the costs). In addition, this simple approach would not be valid for estimating the ICER, which is a non-linear function of the incremental costs and QALYs. For these reasons, we do not consider this approach further. Rather, we present here three flexible strategies for estimating a CACE of the QALYs and the costs, jointly. The first approach combines SURs (equation 1) and 2sls (equation 3) to obtain CACEs for both outcomes accounting for their correlation. This simple approach is known in the econometrics literature as three-stage least squares (3sls).

¶Non-linear versions of the 2sls exist. See for example Clarke and Windmeijer (2012) for an excellent review of methods for binary outcomes.

3.1. Three-stage least squares

Three-stage least squares (3sls) was developed for SUR systems with *endogenous* regressors, i.e. any explanatory variables which are correlated with the error term in equations (1) (Zellner and Theil, 1962). All the parameters appearing in the system are estimated jointly, in three stages. The first two stages are as for 2sls, but with the second stage applied to each of the outcomes.

$$\begin{aligned} \text{(1st stage):} \quad D_i &= \alpha_0 + \alpha_1 Z_i + e_{0i} \\ \text{(2nd stage):} \quad Y_{1i} &= \beta_{01} + \beta_{IV,1} \hat{D}_i + e_{1i} \end{aligned} \tag{4}$$

$$Y_{2i} = \beta_{02} + \beta_{IV,2} \hat{D}_i + e_{2i} \tag{5}$$

As with 2sls, the models can be extended to include baseline covariates. The third stage is the same step used on a SUR with exogenous regressors (equation (1)) for estimating the covariance matrix of the error terms from the two equations (4) and (5). Thus, because we are assuming that Z satisfies the identification assumptions (i)-(v), Z is independent of the residuals at first and second stage, i.e. $Z \perp\!\!\!\perp e_{0i}$, $Z \perp\!\!\!\perp e_{1i}$, and $Z \perp\!\!\!\perp e_{2i}$. Then, the 3sls procedure allows us to obtain the covariance matrix between the residuals e_{1i} and e_{2i} . As with SURs, the 3sls approach does not require to make any distributional assumptions, as estimation can be done by FGLS, and it is robust to heteroscedasticity of the errors in the linear models for the outcomes (Greene, 2002). We note that the 3sls estimator based on FGLS is consistent only if the error terms in each equation of the system and the instrument are independent, which is likely to hold here, as we are dealing with a randomised instrument. In settings where this condition is not satisfied, other estimation approaches such as generalised methods of moments (GMM) warrant consideration (Schmidt, 1990). In the just-identified case, i.e. when there are as many endogenous regressors as there are instruments, classical theory about 3sls estimators shows that the GMM and the FGLS estimators coincide (Greene, 2002). As the 3sls method uses an estimated variance-covariance matrix, it is only asymptotically efficient (Greene, 2002).

3.2. Naive Bayesian estimators

Bayesian models have a natural appeal for cost-effectiveness analyses, as they afford us the flexibility to estimate bivariate models on the expectations of the two outcomes using different distributions, as proposed by (Nixon and Thompson, 2005). These models are often specified by writing a marginal model for one of the outcomes, e.g. the costs Y_1 , and then, a model for Y_2 , conditional on Y_1 .

For simplicity of exposition, we begin by assuming normality for both outcomes and no adjustment for covariates. We have a marginal model for Y_1 and a model for Y_2 conditional on Y_1 (Nixon and Thompson, 2005)

$$Y_{1i} \sim N(\mu_{1i}, \sigma_1^2) \quad \mu_{1i} = \beta_{0,1} + \beta_{1,1} \text{treat}_i \tag{6}$$

$$Y_{2i} | Y_{1i} \sim N(\mu_{2i}, \sigma_2^2(1 - \rho^2)) \quad \mu_{2i} = \beta_{0,2} + \beta_{1,2} \text{treat}_i + \beta_{2,2}(y_{1i} - \mu_{1i}), \tag{7}$$

where ρ is the correlation between the outcomes. The linear relationship between the two outcomes is represented by $\beta_{2,2} = \rho \frac{\sigma_2}{\sigma_1}$.

Because of the non-compliance, in order to obtain a causal estimate of treatment, we need to add a linear model for the treatment received, dependent on randomisation Z_i , similar to the first equation of a 2sls. Formally, this model (denoted uBN, for unadjusted Bayesian Normal) can be written with three equations as follows:

$$\begin{aligned} D_i &\sim N(\mu_{0i}, \sigma_0^2) & \mu_{0i} &= \beta_{0,0} + \beta_{1,0}Z_i \\ Y_{1i} &\sim N(\mu_{1i}, \sigma_1^2) & \mu_{1i} &= \beta_{0,1} + \beta_{1,1}\mu_{0i} \\ Y_{2i} | Y_{1i} &\sim N(\mu_{2i}, \sigma_2^2(1 - \rho^2)) & \mu_{2i} &= \beta_{0,2} + \beta_{1,2}\mu_{0i} + \beta_{2,2}(y_{1i} - \mu_{1i}) \end{aligned} \quad (8)$$

This model is a bivariate version of the ‘unadjusted Bayesian’ method previously proposed for Mendelian randomisation (Burgess and Thompson, 2012). It is called unadjusted, because the variance structure of the outcomes is assumed to be independent of the treatment received. The causal treatment effect for outcome Y_ℓ , with $\ell \in \{1, 2\}$, is represented by $\beta_{1,\ell}$ in equations (8). We use the Fisher’s z-transform of ρ , i.e. $z = \frac{1}{2} \log\left(\frac{1+\rho}{1-\rho}\right)$, for which we assume a vague normal prior, i.e. $z \sim N(0, 10^2)$. We also use vague multivariate normal priors for the regression coefficient (with a precision of 0.01). For standard deviations, we use $\sigma_j \sim \text{Unif}(0, 10)$, for $j \in \{0, 1, 2\}$. This is similar to the priors used in (Lancaster, 2004), and are independent of the regression coefficient of treatment received on treatment allocation $\beta_{1,0}$.

Cost data are notoriously right-skewed, and Gamma-distributions are often used to model them. Thus, we can relax the normality assumption of equation (8), and model Y_1 (i.e. cost) with a Gamma distribution, and treatment received (binary), with a logistic regression. The health outcomes, Y_2 , are still modelled with a normal distribution, as is customary. Because we are using a non-linear model for the treatment received, we use the predicted raw residuals from this model as extra regressors in the outcome models, similar to the 2-stage residual inclusion estimator (Terza et al., 2008). We model Y_1 by its marginal distribution (Gamma) and Y_2 by a conditional Normal distribution, given Y_1 (Nixon and Thompson, 2005). We call this model ‘unadjusted Bayesian Gamma-Normal’ (uBGN) and write it as follows

$$\begin{aligned} D_i &\sim \text{Bern}(\pi_i) & \text{logit}(\pi_i) &= \alpha_0 + \alpha_1 Z_i \\ Y_{1i} &\sim \text{Gamma}(\nu_{1i}, \kappa_1) & r_i &= D_i - \pi_i \\ Y_{2i} | Y_{1i} &\sim N(\mu_{2i}, \sigma_2^2(1 - \rho^2)) & \mu_{1i} &= \beta_{0,1} + \beta_{1,1}D_i + \beta_{1,r}r_i \\ & & \mu_{2i} &= \beta_{0,2} + \beta_{1,2}D_i + \beta_{2,r}r_i + \beta_{2,2}(y_{1i} - \mu_{1i}) \end{aligned} \quad (9)$$

where $\mu_1 = \frac{\nu_1}{\kappa_1}$, is the mean of the Gamma distributed costs, with shape ν_1 and rate κ_1 . Again, we express $\beta_{2,2} = \rho \frac{\sigma_2}{\sigma_1}$, and assume a vague Normal prior on the Fisher’s z-transform of ρ , $z \sim N(0, 10^2)$. The prior distribution for ν_1 is $\text{Gamma}(0.01, 0.01)$. We also assume a Gamma prior for the intercept term of the cost equation, $\beta_{0,1} \sim \text{Gamma}(0.01, 0.01)$. All the other regression parameters have the same priors as those used in the uBN model.

The models introduced in this section, uBN and uBGN, are estimated in one stage, allowing feedback between the regression equations and the propagation of uncertainty. However, these 'unadjusted' methods ignore the correlation between the outcomes and the treatment received. This misspecification of the covariance structure may result in biases in the causal effect, which are likely to be more important at higher levels of non-compliance.

3.3. Bayesian simultaneous equations (BFL)

We now introduce an approach that models the covariance between treatment received and outcomes appropriately, using a system of simultaneous equations. This can be done via full or limited information maximum likelihood, or by using MCMC to estimate the parameters in the model simultaneously allowing for proper Bayesian feedback and propagation of uncertainty. Here, we propose a Bayesian approach which is an extension of the methods presented in Burgess and Thompson (2012); Kleibergen and Zivot (2003); Lancaster (2004).

This method treats the endogenous variable D and the cost-effectiveness outcomes as covariant and estimates the effect of treatment allocation, as follows. Let $(D_i, Y_{1i}, Y_{2i})^\top$ be the transpose of vector of outcomes, which now includes treatment received, as well as the bivariate endpoints of interest. We treat all three variables as multivariate normally distributed, so that

$$\begin{pmatrix} D_i \\ Y_{1i} \\ Y_{2i} \end{pmatrix} \sim N \left\{ \begin{pmatrix} \mu_{0i} \\ \mu_{1i} \\ \mu_{2i} \end{pmatrix}, \Sigma = \begin{pmatrix} \sigma_0^2 & s_{01} & s_{02} \\ s_{01} & \sigma_1^2 & s_{12} \\ s_{02} & s_{12} & \sigma_2^2 \end{pmatrix} \right\}; \quad \begin{aligned} \mu_{0i} &= \beta_{0,0} + \beta_{1,0}Z_i \\ \mu_{1i} &= \beta_{0,1} + \beta_{1,1}\beta_{1,0}Z_i \\ \mu_{2i} &= \beta_{0,2} + \beta_{1,2}\beta_{1,0}Z_i \end{aligned} \quad (10)$$

where $s_{ij} = \text{cov}(Y_i, Y_j)$, and the causal treatment effect estimates are $\beta_{1,1}$ and $\beta_{1,2}$ respectively. For the implementation, we use vague normal priors for the regression coefficients, i.e. $\beta_{m,j} \sim N(0, 10^2)$, for $j \in \{0, 1, 2\}$, $m \in \{0, 1\}$, and a Wishart prior for the inverse of Σ (Gelman and Hill, 2006).

4. Simulation study

We now use a factorial simulation study to assess the finite sample performance of the alternative methods. The first factor is the proportion of participants who do not comply with the experimental regime, when assigned to it, expressed as a percentage of the total (one-sided non-compliance). Bias is expected to increase with increasing levels of non-compliance. A systematic review (Dodd et al., 2012) found that the percentage of non-compliance was less than 30% in two-thirds of published RCTs, but greater than 50% in one-tenth of studies. Here, two levels of non-compliance are chosen, 30% and 70%. As costs are typically skewed, three different distributions (Normal, Gamma or Inverse Gaussian – IG) are used to simulate cost data. As the 2sls approach fails to accommodate the correlation between the endpoints, we examined the impact of different

levels of correlation on the methods' performance; ρ takes one of four values ± 0.4 , ± 0.8 . The final factor is the sample size of the RCT, taking two settings $n = 100$, and 1000. In total, there are $2 \times 3 \times 4 \times 2 = 48$ simulated scenarios.

To generate the data, we begin by simulating $U \sim N(0.50, 0.25^2)$, independently from treatment allocation. U represents a pre-randomisation variable that is a common cause of both the outcomes and the probability of non-compliance, *i.e.* it is a confounding variable, which we assume is unobserved.

Now, let $S_i \sim \text{Bern}(\pi_s)$ be the random variable denoting whether the i th individual switches from allocated active treatment to control. The probability π_s of one-way non-compliance with allocated treatment depends on U , in the following way,

$$\pi_s = \begin{cases} p + 0.1, & \text{if } u > 0.5, \\ p - 0.1, & \text{otherwise} \end{cases} \quad (11)$$

where p denotes the corresponding average non-compliance percentage expressed as a probability, *i.e.* here $p \in \{0.3, 0.7\}$. We now generate D_i , the random variable of treatment received as

$$D_i = \begin{cases} Z_i, & \text{if either } s_i = 0 \text{ or } Z_i = 0, \\ 1 - Z_i, & \text{if } s_i = 1 \text{ and } Z_i = 1 \end{cases} \quad (12)$$

where Z_i denotes the random allocation for subject i .

Then, the means for both outcomes are assumed to depend linearly on treatment received and the unobserved confounder U as follows:

$$\mu_1 = E[Y_1] = 1.2 + 0.4D_i + 0.16(u_i - 0.5) \quad (13)$$

$$\mu_2 = E[Y_2] = 0.5 + 0.2D_i + 0.04(u_i - 0.5) \quad (14)$$

Finally, the bivariate outcomes are generated using Gaussian copulas, initially with normal marginals. In subsequent scenarios, we consider Gamma or Inverse Gaussian marginals for Y_1 and normal for Y_2 . The conditional correlation between the outcomes, ρ , is set according to the corresponding scenario.

For the scenarios where the endpoints are assumed to follow a bivariate normal distribution, the variances of the outcomes are set to $\sigma_1^2 = 0.2^2$, $\sigma_2^2 = 0.1^2$ respectively, while for scenarios with Gamma and IG distributed Y_1 , the shape parameter is $\eta = 4$. For the Gamma case, this gives a variance for Y_1 equal to $\sigma_1^2 = 0.36$ in control and $\sigma_1^2 = 0.64$ in the intervention group. When $Y_1 \sim \text{IG}(\mu_1, \eta)$, the expected variance in the control group is $\sigma_1^2 = 0.432$, and $\sigma_1^2 = 1.024$ in those receiving the intervention.

The simulated endpoints represent cost-effectiveness variables that have been rescaled for computational purposes, with costs divided by 1000, and QALYs by 0.1, such that the true values are £400 (incremental costs), and 0.02 (incremental QALYs) and so with a threshold value of $\lambda = \text{£}30000$ per QALY, the true causal INB is £200.

For each simulated scenario, we obtained $M = 2500$ sets. For the Bayesian analyses, we use the median of the posterior distribution as the 'estimate' of the parameter of interest,

and the standard deviation of the posterior distribution as the standard error. Equal tailed 95% posterior credible intervals are also obtained. We use the term confidence interval for the Bayesian credible intervals henceforth, to have a unified terminology for both Bayesian and frequentist intervals.

Once the corresponding causal estimate has been obtained in each of the 2500 replicated sets under each scenario in turn, we compute the median bias of the estimates, coverage of 95% confidence intervals (CI), median CI width and root mean square error (RMSE). We report median bias as opposed to mean bias, because the BFL leads to a posterior distribution of the causal parameters which is Cauchy-like (Kleibergen and Zivot, 2003). A method is ‘adequate’, if it results in low levels of bias (median bias $\leq 5\%$) with coverage rates within 2.5% of the nominal value.

Implementation:

The 3sls was fitted using `systemfit` package in R using FGLS, and the Bayesian methods were run using JAGS from R (`r2jags`). Two chains, each one with 5000 initial iterations and 1000 burn-in were used. The multiple chains allowed for a check of convergence by the degree of their mixing and the initial iterations enabled to estimate iteration autocorrelation. A variable number of further 1000-iteration runs were performed until convergence was reached as estimated by the absolute value of the Geweke statistics for the first 10% and last 50% of iterations in a run being below 2.5. A final additional run of 5000 iterations was performed for each chain to achieve a total sample of 10000 iterations, and a MC error of about 1% of the parameter SE on which to base the posterior estimates. For the uBGN, an offset of 0.01 is added to the shape parameter ν_1 for the Gamma distribution of the cost, to prevent the sampled shape parameter to become too close to zero, which may result in infinite densities. See the Supplementary File for the JAGS model code for BFL.

4.1. Simulation Study Results

Bias

Figure 1 shows the median bias corresponding to scenarios with 30% non-compliance, by cost distributions (left to right) and levels of correlation between the two outcomes, for sample sizes of $n = 100$ (upper panel) and $n = 1000$ (lower panel). With the larger sample size, for all methods, bias is negligible with normally distributed costs, and remain less than 5% when costs are Gamma-distributed. However, when costs follow an Inverse Gaussian distribution, and the absolute levels of correlation between the endpoints are high (± 0.8), the uBGN approach results in biased estimates, around 10% bias for the estimated incremental cost, and between 20 and 40% for the estimated INB. With the small sample size and when costs follow a Gamma or Inverse Gaussian distribution, both unadjusted Bayesian methods provide estimates with moderate levels of bias. With 70% non-compliance (Figure A3 in the Supplementary file), the unadjusted methods result in important biases which persist even with large sample sizes, especially for scenarios with non-normal outcomes. For small sample settings, uBN reports positive bias (10 to 20%)

in the estimation of incremental QALYs, and the resulting INB, irrespective of the cost distribution. The uBGN method reports relatively unbiased estimates of the QALYs, but large positive bias (up to 60%) in the estimation of costs, and hence, there is substantial bias in the estimated INB (up to 200%). The unadjusted Bayesian methods ignore the positive correlation between the confounding variable and both the treatment received and the outcome. These methods therefore provide estimates of the casual effects that exceed the true values, i.e. have a positive bias. By contrast, the BFL and the 3sls provide estimates with low levels of bias across most settings.

CI coverage and width

Table 2 presents the results for CI coverage and width, for scenarios with a sample size of $n = 100$, absolute levels of correlation between the endpoints of 0.4, and 30% non-compliance. All other results are shown in the Supplementary file. The 2sls INB ignores the correlation between costs and QALYs, and thus, depending on the direction of this correlation, 2sls reports CI coverage that is above (positive correlation) or below (negative correlation) nominal levels. This divergence from nominal levels increases with higher absolute levels of correlation (see Supplementary file, Table A6).

The uBN approach results in over-coverage across many settings, with wide CIs. For example, for both levels of non-compliance and either sample size, when the costs are Normal, the CI coverage rates for both incremental costs and QALYs exceed 0.98. The interpretation offered by Burgess and Thompson (2012) is also relevant here: the uBN assumes that the treatment received and the outcomes variance structures are uncorrelated, and so when the true correlation is positive, the model overstates the variance and leads to wide CIs. By contrast, the uBGN method results in low CI coverage rates for the estimation of incremental costs, when costs follow an inverse Gaussian distribution. This is because the model incorrectly assumes a Gamma distribution, thereby underestimating the variance. The extent of the under-coverage appears to increase with higher absolute values of the correlation between the endpoints, with coverage as low as 0.68 (incremental costs) and 0.72 (INB) in scenarios where the absolute value of correlation between costs and QALYs is 0.8. (see Supplementary file, Table A7).

The BFL approach reports estimates with CI coverage close to the nominal when the sample size is large, but with excess coverage (greater than 0.975), and relatively wide CI, when the sample size is $n = 100$ (see Table 2 for 30% non-compliance, and Table A4 in the Supplementary file for the result corresponding to 70% non-compliance). By contrast, the 3sls reports CI coverage within 2.5% of nominal levels for each sample size, level of non-compliance, cost distribution and level of correlation between costs and QALYs.

RMSE

Table 3 reports RMSE corresponding to 30% non-compliance, and $n = 100$. The least squares approaches result in lower RMSE than the other methods for the summary statistic of interest, the INB. This pattern is repeated across other settings, see the Supplementary file, Tables A10–A16. ||

||The RMSE for the 2sls and 3sls estimates is the same for each of the outcomes considered,

5. Results for the motivating example

We now compare the methods in practice by applying them to the REFLUX dataset. Only 48% of the individuals have completely observed cost-effectiveness outcomes: there were 185 individuals with missing QALYs, 166 with missing costs, and a further 13 with missing EQ5D₀ at baseline, with about a third of those with missing outcomes having switched from their allocated treatment. These missing data not only bring more uncertainty to our analysis, but more importantly, unless the missing data are handled appropriately can lead to biased causal estimates (Daniel et al., 2012). A complete case analysis would be unbiased, albeit inefficient, if the missingness is conditionally independent of the outcomes given the covariates in the model (White et al., 2010), even when the covariates have missing data, as is the case here.** Alternatively, a more plausible assumption is to assume the missing data are missing at random (MAR), i.e. the probability of missingness depends only on the observed data, and use multiple imputation (MI) or a full Bayesian analysis to obtain valid inferences.

Therefore, we perform MI prior to carrying out 2sls and 3sls analyses. We begin by investigating all the possible associations between the covariates available in the data set and the missingness, univariately for costs, QALYs and baseline EQ5D₀. Covariates which are predictive of both, the missing values and the probability of missing, are to be included in the imputation model as auxiliary variables, as conditioning on more variables helps make the MAR assumption more plausible. None of the available covariates satisfies these criteria and therefore, we do not include any auxiliary variables in our imputation models. Thus, we impute total cost, total QALYs and baseline EQ5D₀, 50 times by chained equations, using predictive mean matching (PMM), taking the 5 nearest neighbours as donors (White et al., 2011), including treatment received in the imputation model and stratifying by treatment allocation. We perform 2sls on costs and QALYs independently and calculate (within MI) SE for the INB assuming independence between costs and QALYs. For the 3sls approach, the model is fitted to both outcomes simultaneously, and the post-estimation facilities are used to extract the variance-covariance estimate, and compute the estimated INB and its corresponding SE. We also use the CACE estimates of incremental cost and QALYs to obtain the ICER. After applying each method to the 50 MI sets, we combine the results using Rubin's rules (Rubin, 1987).††

For the Bayesian approaches, the missing values become extra parameters to model. because the two methods obtain the same point estimate, and hence, by definition, they have the same empirical standard-error, even though they have different model-based standard errors for INB. This is in contrast to the differences observed in the performance of measures based on the CI. Coverage rate and CI width corresponding to these two methods are different for the INB, because the confidence intervals are constructed using the model-based SE. See the Supplementary File for further details.

**This mechanism is a special case of missing not at random.

††Applying IV 2sls and 3sls with multiply imputed datasets, and combining the results using Rubin's rules can be done automatically in Stata using `mi estimate, cmdok: ivregress 2sls` and `mi estimate, cmdok: reg3`. In R, `ivregress` can be used with `with.mids` command

Since baseline EQ5D₀ has missing observations, a model for its distribution is added EQ5D₀ $\sim N(\mu_{q0}, \sigma_{q0}^2)$, with a vaguely informative prior for $\mu_{q0} \sim \text{Unif}(-0.5, 1)$, and an uninformative prior for $|\sigma_{q0}| \sim N(0, 0.01)$. We add two extra lines of code to the models to obtain posterior distributions for INB and ICERs. We center the outcomes around the empirical mean (except for costs, when modelled as Gamma) and re-scale the costs (dividing by 1000) to improve mixing and convergence. We use two chains, initially running 15,000 iterations with 5,000 as burn-in. After checking visually for auto-correlation, an extra 10,000 iterations are needed to ensure that the density plots of the parameters corresponding to the two chains are very similar, denoting convergence to the stationary distribution. Enough iterations for each chain are kept to make the total effective sample (after accounting for auto-correlation) equal to 10,000. ‡‡

Table 1 shows the results for incremental costs, QALYs and INB for the motivating example adjusted for baseline EQ5D₀. Bayesian posterior distribution are summarised by their median value and 95% credible intervals. The CACEs are similar across the methods, except for uBGN, where the incremental QALYs CACE is nearly halved, resulting in a smaller INB with a CI that includes 0. In line with the simulation results, this would suggest that, where the uBGN is misspecified according to the assumed cost distribution, it can provide a biased estimate of the incremental QALYs.

Comparing the CACEs to the ITT, we see that the incremental cost estimates increases between an ITT and a CACE, as actual receipt of surgery carries with it higher costs than the mere offering of surgery does not. Similarly, the incremental QALYs are larger, meaning that amongst compliers, those receiving surgery have a greater gain in quality of life, over the follow-up period. The CACE for costs are relatively close to the per-protocol incremental costs reported in the original study, £2324 (1780, 2848). In contrast, the incremental QALYs according to PP on complete-cases originally reported was 0.3200 (0.0837, 0.5562), considerably smaller than our CACE estimates (Grant et al., 2013). The ITT ICER obtained after MI was £4135, while using causal incremental costs and QALYs, the corresponding estimates of the CACE for ICER were £4140 (3sls), £5189 (uBN), £5960 (uBGN), and £3948 (BFL). The originally reported per-protocol ICER is £7932 per extra QALY was obtained on complete cases only (Grant et al., 2013).

These results may be sensitive to the modelling of the missing data. As a sensitivity analysis to the MAR assumption, we present the complete case analysis in Table A1 in the Supplementary File. The conclusions from complete-case analysis are similar to those obtained under MAR.

We also explore the sensitivity to choices of priors, by re-running the BFL analyses using different priors, first for the multivariate precision matrix, keeping the priors for

within mice, but `systemfit` cannot presently be combined with this command so, Rubin's rules have to be coded manually. Sample code is available in the Supplementary File.

‡‡Multivariate normal nodes cannot be partially observed in JAGS, thus, we run BFL models on all available data within WinBUGs. An observation with zero costs was set to missing when running the Bayesian Gamma model, which requires strictly positive costs.

the coefficients normal, and then a second analysis, with uniform priors for the regression coefficient, and an inverse Wishart prior with 6 degrees of freedom and a identity scale matrix, for the precision. The results are not materially changed (see Supplementary File, Table A2).

The results of the within-trial CEA suggest that amongst compliers, laparoscopy is more cost-effective than medical management for patients suffering from GORD. The results are robust to the choice of priors, and to the assumptions about the missing data mechanism. The results for the uBGN differ somewhat from the other models, and as our simulations show, the concern is that such unadjusted Bayesian models are prone to bias from model misspecification.

6. Discussion

This paper extends existing methods for CEA (Willan et al., 2003; Nixon and Thompson, 2005), by providing IV approaches for obtaining causal cost-effectiveness estimates for RCTs with non-compliance. The methods developed here however are applicable to other settings with multivariate continuous outcomes more generally, for example RCTs in education, with different measures of attainment being combined into an overall score. To help dissemination, we provide code in the Supplementary File.

We proposed exploiting existing 3sls methods and also considered IV Bayesian models, which are extensions of previously proposed approaches for univariate continuous outcomes. Burgess and Thompson (2012) found the BFL was median unbiased and gave CI coverage close to nominal levels, albeit with wider CIs than least-squares methods. Their ‘unadjusted Bayesian’ method, similar to our uBN approach, assumes that the error term for the model of treatment received on treatment allocated is uncorrelated with the error from the outcome models. This results in bias and affects the CI coverage. Our simulation study shows that, in a setting with multivariate outcomes, the bias can be substantial. A potential solution to this could be using priors for the error terms that reflect the dependency of the error terms explicitly. For example, Rossi et al. (2012) propose using a prior for the errors that explicitly depends on the coefficient $\beta_{1,0}$, the effect of treatment allocation on treatment received, in equation (8). Kleibergen and Zivot (2003) propose priors that also reflect this dependency explicitly, and replicate better the properties of the 2sls. This is known as the “Bayesian two-stage approach”.

The results of our simulations show that applying 2sls separately to the univariate outcomes leads to inaccurate 95% CI around the INB, even with moderate levels of correlation between costs and outcomes (± 0.4). Across all the settings considered, the 3sls approach resulted in low levels of bias for the INB and unlike 2sls, provided CI coverage close to nominal levels. BFL performed well with large sample sizes, but produced standard deviations which were too large when the sample size was small, as can be seen from the over-coverage, with wide CIs.

The REFLUX study illustrated a common concern in CEA, in that the levels of non-

compliance in the RCT were different, in this case higher, to those in routine practice. The CACEs presented provide the policy-maker with an estimate of what the relative cost-effectiveness would be if all the RCT participants had complied with their assigned treatment, which is complementary to the ITT estimate. Since we judged the IV assumptions for identification likely to hold in this case-study, we conclude that either 3sls or BFL provide valid inferences for the CACE of INB. The re-analysis of the REFLUX case study also provided the opportunity to investigate the sensitivity to the choice of priors in practice. Here we found that our choice of weakly informative priors, which were relatively flat in the region where the values of the parameters were anticipated to be, together with samples of at least size 100, had minimal influence on the posterior estimates. We repeated the analysis using different vague priors for the parameters of interests and the corresponding results were not materially changed.

The REFLUX study also illustrated a further complication that may arise in practice, namely that covariate or outcome data are missing. Here we illustrated how the methods for estimating the CACE can also accommodate missing data, under the assumption that the data are missing at random (MAR), without including any auxiliary variables in the imputation or Bayesian models. However, more generally, where there are auxiliary variables available, these should be done included in the imputation or Bayesian models. If the auxiliary variables have missing values themselves, this can be accommodated easily via chained-equations MI, but for the Bayesian approach, an extra model for the distribution of the auxiliary variable, given the other variables in the substantive model and the outcome needs to be added.

We considered here relatively simple frequentist IV methods, namely 2sls and 3sls. One alternative approach to the estimation of CACEs for multivariate responses, is to use linear structural equation modelling, estimated by maximum-likelihood expectation-maximization (ML-EM) algorithm (Jo and Muthén, 2001). Further, we only considered those settings where a linear additive treatment effect is of interest, and the assumptions for identification are met. Where interest lies in systems of simultaneous non-linear equations with endogeneous regressors, GMM or generalised structural equation models can be used to estimate CACEs (Davidson and MacKinnon, 2004).

There are several options to study the sensitivity to departures from the identification assumptions. For example, if the exclusion restriction does not hold, a Bayesian parametric model can use priors on the non-zero direct effect of randomisation on the outcome for identification (Conley et al., 2012; Hirano et al., 2000). Since the models are only weakly identified, the results would depend strongly on the parametric choices for the likelihood and the prior distributions. In the frequentist IV framework, such modelling is also possible, see Baiocchi et al. (2014) for an excellent tutorial on how to conduct sensitivity analysis to violations of the ER and monotonicity assumptions. Alternatively, violations of the ER can also be handled by using baseline covariates to model the probability of compliance directly, within structural equation modelling via ML-EM framework (Jo, 2002a,b).

Settings where the instrument is only weakly correlated with the endogenous variable have not been considered here, because for binary non-compliance with binary allocation, the percentage of one-way non-compliance would need to be in excess of 85%, for the F-statistic of the randomisation instrument to be less than 10, the traditional cutoff beneath which an instrument is regarded as ‘weak’. Such levels of non-compliance are not realistic in practice, with the reported median non-compliance equal to 12% (Zhang et al., 2014). Nevertheless, Bayesian IV methods have been shown to perform better than 2sls methods when the instrument is weak (Burgess and Thompson, 2012).

Also, for simplicity, we restricted our analysis of the case study to MAR and complete cases assumptions. Sensitivity to departures from these assumptions is beyond the scope of this paper, but researchers should be aware of the need to think carefully about the possible causes of missingness, and conduct sensitivity analysis under MNAR, assuming plausible differences in the distributions of the observed and the missing data. When addressing the missing data through Bayesian methods, the posterior distribution can be sensitive to the choice of prior distribution, especially with a large amount of missing data (Hirano et al., 2000).

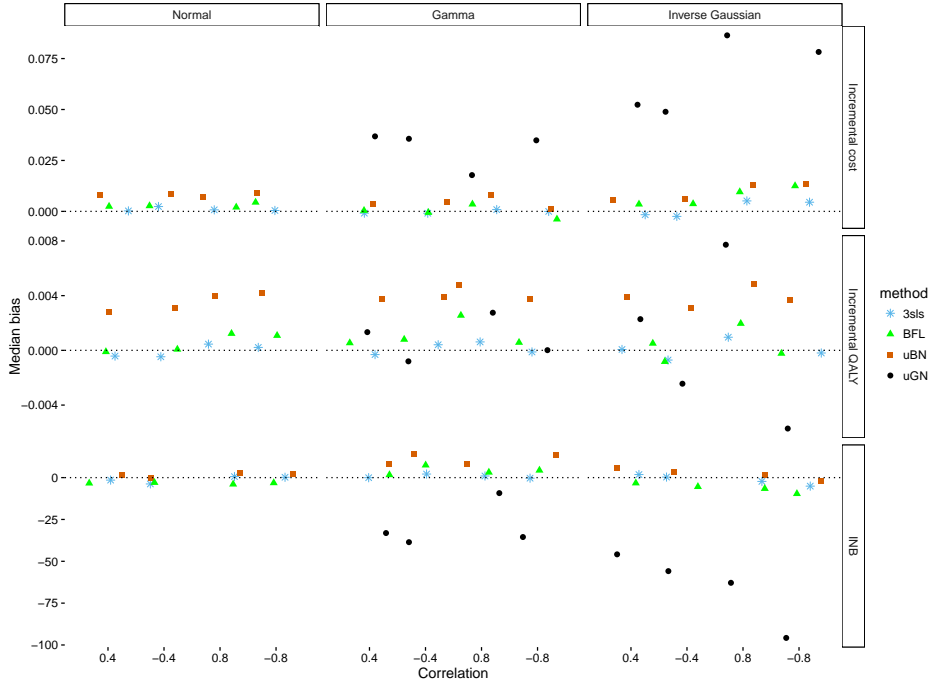
Future research directions could include exploiting the additional flexibility of the Bayesian framework to incorporate informative priors, perhaps as part of a comprehensive decision modelling approach. The methods developed here could also be extended to handle time-varying non-compliance.

Acknowledgements

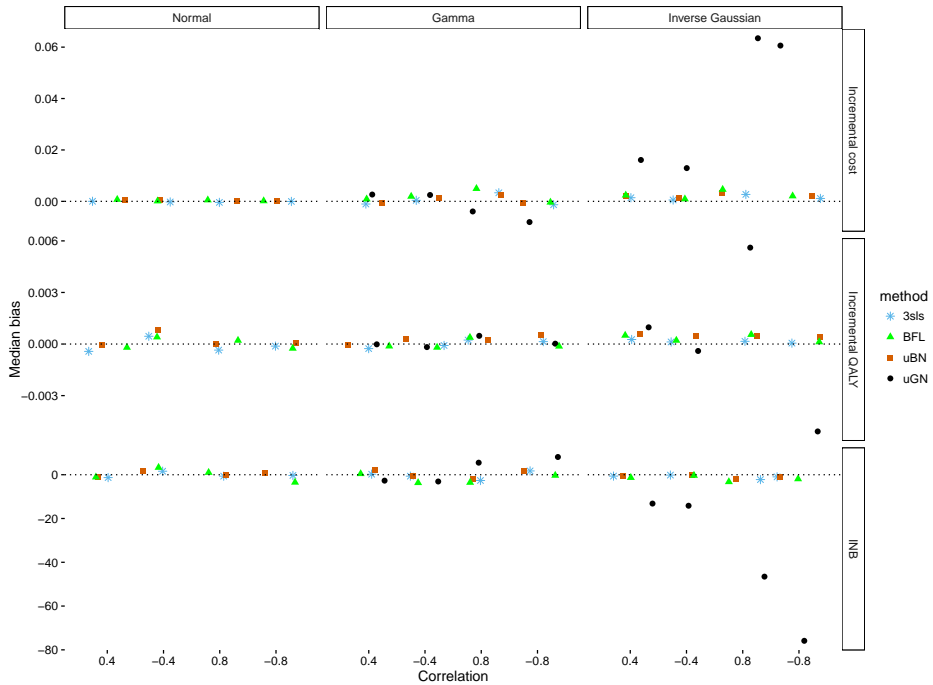
We thank Mark Sculpher, Rita Faria, David Epstein, Craig Ramsey and the REFLUX study team for access to the data. We also thank James Carpenter and Simon Thompson for comments on earlier drafts.

Karla DiazOrdaz was supported by UK Medical Research Council Career development award in Biostatistics MR/L011964/1. This report is independent research supported by the National Institute for Health Research (Senior Research Fellowship, Richard Grieve, SRF-2013-06-016). The views expressed in this publication are those of the author(s) and not necessarily those of the NHS, the National Institute for Health Research or the Department of Health.

Fig. 1. Median Bias for scenarios with 30% non-compliance and sample sizes of (a) $n = 100$ (top) and (b) $n = 1000$. Results are stratified by cost distribution, and correlation between cost and QALYs. The dotted line represents zero bias. Results for 2sls (not plotted) are identical to those for 3sls; uBGN was not applied to Normal cost data.



(a) $n=100$



(b) $n=1000$

Table 1. The REFLUX study: descriptive statistics and cost-effectiveness according to ITT and alternative methods for estimating the CACE. Follow-up period is five years, and treatment switches are defined within the first year post randomisation. Costs and INB numbers rounded to the nearest integer.^a

	Medical management	Laparoscopic surgery
N Assigned	179	178
N (%) Switched	10 (8.3)	67 (28.4)
N (%) missing costs	83 (46)	83 (47)
Mean (SD) observed cost in £	1258 (1687)	2971 (1828)
N (%) missing QALYs	91 (51)	94 (53)
Mean (SD) observed QALYs	3.52 (0.99)	3.74 (0.90)
<i>Baseline variables</i>		
N (%) missing EQ5D ₀	6 (3)	7 (4)
Mean (SD) observed EQ5D ₀	0.72 (0.25)	0.71 (0.26)
Correlation between costs and QALYs	-0.42	-0.07
Correlation of costs and QALYs by treatment received	-0.36	-0.18
Incremental costs, QALYs and INB of surgery vs medicine		
Outcome	Method	estimate (95% CI)
Incremental cost	ITT	1103 (593, 1613)
	2sls	1899 (1073, 2724)
	3sls	1899 (1073, 2724)
	uBN	2960 (2026, 3998)
	uBGN	2176 (1356, 3031)
	BFL	2030 (1170, 2878)
	Incremental QALYs	ITT
2sls		0.516 (0.103, 0.929)
3sls		0.516 (0.103, 0.929)
uBN		0.568 (0.181, 0.971)
uBGN		0.268 (-0.229, 0.759)
BFL		0.511 (0.121, 0.947)
INB		ITT
	2sls	13587 (1101, 26073)
	3sls	13587 (1002, 26173)
	uBN	14091 (2485, 26086)
	uBGN	5869 (-9204, 20740)
	BFL	13340 (1406, 26315)

^a uBN: unadjusted Bayesian Normal-Normal model; uBGN: unadjusted Bayesian Gamma-Normal models; BFL: Bayesian Full likelihood models.

Table 2. CI Coverage rates (CR) and median width for incremental cost, QALYs, and INB, across scenarios with 30% non-compliance, sample size $n = 100$ and moderate correlation ρ between outcomes and even rows to negative). uBGN was not applied in settings with normal cost data.^a

	ρ	2sls		3sls		uBN		uBGN		BFL	
		CR	CIW	CR	CIW	CR	CIW	CR	CIW	CR	CIW
$Y_1 \sim N$											
Cost	0.4	.952	.228	.952	.228	.992	.312			.988	.299
	-0.4	.952	.229	.952	.229	.993	.325			.986	.297
QALYs	0.4	.946	.112	.946	.112	.988	.155			.950	.121
	-0.4	.950	.113	.950	.113	.992	.163			.950	.121
INB	0.4	.988	405	.953	319	.982	398			.966	376
	-0.4	.900	409	.948	475	.951	509			.962	525
$Y_1 \sim G$											
Cost	0.4	.952	.756	.952	.756	.955	.815	.941	.818	.954	.823
	-0.4	.942	.759	.942	.759	.949	.828	.936	.822	.945	.811
QALYs	0.4	.959	.113	.959	.113	.993	.160	.960	.122	.960	.122
	-0.4	.959	.113	.949	.113	.995	.163	.954	.122	.954	.122
INB	0.4	.982	829	.948	696	.958	764	.942	748	.956	760
	-0.4	.914	833	.948	943	.930	921	.941	1019	.951	1014
$Y_1 \sim IG$											
Cost	0.4	.951	.880	.951	.880	.958	.949	.904	.866	.956	.945
	-0.4	.950	.878	.950	.878	.958	.951	.905	.864	.954	.932
QALYs	0.4	.945	.112	.945	.112	.991	.161	.944	.120	.999	.206
	-0.4	.954	.112	.954	.112	.993	.161	.952	.120	.999	.204
INB	0.4	.980	944	.954	818	.959	889	.917	814	.984	1001
	-0.4	.917	942	.947	1049	.934	1034	.911	1041	.971	1203

^a uBN: unadjusted Bayesian Normal-Normal model, uBGN: unadjusted Bayesian Gamma-Normal models;

BFL: Bayesian Full likelihood models.

Table 3. RMSE for incremental Cost, QALYs and INB across scenarios with 30% non-compliance, moderate correlation between outcomes and sample size $n = 100$. uBGN was not applied in settings with normal cost data. Numbers for INB have been rounded to the nearest integer.^a

Cost distribution	ρ	3sls ^b	uBN	uBGN	BFL
Normal	Cost				
	0.4	0.058	0.060		0.059
	-0.4	0.060	0.062		0.061
	QALYs				
	0.4	0.029	0.030		0.030
	-0.4	0.029	0.030		0.030
	INB				
	0.4	83	84		87
	-0.4	125	127		125
	Gamma	Cost			
0.4		0.198	0.202	0.212	0.202
-0.4		0.200	0.204	0.212	0.203
QALYs					
0.4		0.030	0.030	0.030	0.029
-0.4		0.029	0.030	0.030	0.030
INB					
0.4		181	184	193	184
-0.4		246	251	261	252
Inverse Gaussian		Cost			
	0.4	0.230	0.232	0.252	0.232
	-0.4	0.230	0.232	0.250	0.232
	QALYs				
	0.4	0.029	0.030	0.030	0.030
	-0.4	0.029	0.030	0.030	0.030
	INB				
	0.4	211	214	231	214
	-0.4	273	278	296	278

^a uBN: unadjusted Bayesian Normal-Normal model; uBGN: unadjusted Bayesian Gamma-Normal models; BFL: Bayesian Full Likelihood

^b The RMSE corresponding to 2sls is identical to that for 3sls, by definition.

References

- Angrist, J. D., Imbens, G. W. and Rubin, D. B. (1996), ‘Identification of causal effects using instrumental variables’, *Journal of the American Statistical Association* **91**(434), pp. 444–455.
- Angrist, J. D., Pischke, J. S. (2008), *Mostly Harmless Econometrics: An Empiricist’s Companion*. Princeton University Press.
- Baiocchi, M., Cheng, J. and Small, D. S. (2014), ‘Instrumental variable methods for causal inference’, *Statistics in Medicine* **33**(13), 2297–2340.
- Brilleman, S., Metcalfe, C., Peters, T. and Hollingsworth, W. (2015), ‘The reporting of treatment non-adherence and its associated impact on economic evaluations conducted alongside randomised trials: a systematic review.’, *Value in Health* **19** (1), pp., 99 – 108.
- Burgess, S. and Thompson, S. G. (2012), ‘Improving bias and coverage in instrumental variable analysis with weak instruments for continuous and binary outcomes’, *Statistics in Medicine* **31**(15), 1582–1600.
- Clarke, P. S. and Windmeijer, F. (2012), ‘Instrumental Variable Estimators for Binary Outcomes.’ *Journal of the American Statistical Association* 107(500):1638–1652.
- Conley, T. G., Hansen, C. B. and Rossi, P. E. (2012), ‘Plausibly exogenous’, *Review of Economics and Statistics* **94**(1), 260–272.
- Daniel, R. M., Kenward, M. G., Cousens, S. N. and De Stavola, B. L. (2012), ‘Using causal diagrams to guide analysis in missing data problems’, *Statistical Methods in Medical Research*: 21(3), 243–256.
- Davidson, R. and MacKinnon, J. G. (2004). *Economic theory and methods*, New York: Oxford University Press.
- Didelez, V., Meng, S. and Sheehan, N. (2010) ‘ Assumptions of IV Methods for Observational Epidemiology’. *Statistical Science*, **25**(1), 22–40.
- Dodd, S., White, I. and Williamson, P. (2012), ‘Nonadherence to treatment protocol in published randomised controlled trials: a review’, *Trials* **13**(1), 84.
- Gelman, A. and Hill, J. (2006), *Data Analysis Using Regression and Multilevel/Hierarchical Models*, Analytical Methods for Social Research, Cambridge University Press.
- Grant, A. M., Boachie, C., Cotton, S. C., Faria, R., Bojke, L. and Epstein, D. (2013) ‘Clinical and economic evaluation of laparoscopic surgery compared with medical management for gastro-oesophageal reflux disease: a 5-year follow-up of multicentre randomised trial (the REFLUX trial).’, *Health Technology Assessment* **17**(22).
- Grant, A., Wileman, S., Ramsay, C., Boyke, L., Epstein, D. and Sculpher, M. (2008), ‘The effectiveness and cost-effectiveness of minimal access surgery amongst people with gastro-oesophageal reflux disease- a uk collaborative study. the REFLUX trial.’ *Health Technology Assessment* **12** (31).

- Greene, W. (2002). *Econometric Analysis*, Prentice-Hall international editions, Prentice Hall.
- Hernán, M. A. and Robins, J. M. ‘Instruments for causal inference: an epidemiologist’s dream?’, *Epidemiology* **17** (4), 360–372.
- Hirano, K. , Imbens, G. W. , Rubin, D. B. and Zhou, X. H. ‘Assessing the effect of an influenza vaccine in an encouragement design’, *Biostatistics* **1**, 69 –88.
- Hoch, J. S., Briggs, A. H. and Willan, A. R., (2002). ‘Something old, something new, something borrowed, something blue: A framework for the marriage of health econometrics and costeffectiveness analysis’. *Health economics*, **11** (5) 415–430.
- Hughes, D., Charles, J., Dawoud, D., Edwards, R. T., Holmes, E. , Jones, C. , Parham, P. , Plumpton, C. , Ridyard, C., Lloyd-Williams, H., Wood, E., and Yeo, S. T. (2016). ‘Conducting economic evaluations alongside randomised trials: Current methodological issues and novel approaches.’ *PharmacoEconomics*, 1–15.
- Imbens, G. W. and Angrist, J. D. (1994), ‘Identification and estimation of local average treatment effects’, *Econometrica* **62**(2), 467–475.
- Imbens, G. W. and Rubin, D. B. (1997), ‘Bayesian inference for causal effects in randomized experiments with noncompliance’, *The Annals of Statistics* **25**(1), 305–327.
- Jo B. (2002a) ‘Estimating intervention effects with noncompliance: Alternative model specifications’. *Journal of Educational and Behavioral Statistics*. **27**:385 –420.
- Jo B. (2002b) ‘Model misspecification sensitivity analysis in estimating causal effects of interventions with noncompliance.’ *Statistics in Medicine*. **21**: 3161– 3181.
- Jo B. and Muthén B. O. (2001) ‘Modeling of intervention effects with noncompliance: a latent variable approach for randomised trials,’ *In Marcoulides GA, Schumacker RE, eds. New developments and techniques in structural equation modeling*. Lawrence Erlbaum Associates, Mahwah, New Jersey: 57–87.
- Kleibergen, F. and Zivot, E. (2003), ‘Bayesian and classical approaches to instrumental variable regression’, *Journal of Econometrics* **114**(1), 29 – 72.
- Lancaster, T. (2004), *Introduction to Modern Bayesian Econometrics*, Wiley.
- Latimer, N. R., Abrams, K., Lambert, P., Crowther, M., Wailoo, A., Morden, J., Akehurst, R. and Campbell, M. (2014), ‘Adjusting for treatment switching in randomised controlled trials - a simulation study and a simplified two-stage method’, *Statistical Methods in Medical Research*: 0962280214557578.
- Mantopoulos, T., Mitchell, P.M., Welton, N.J. McManus, R. and Andronis, L. (2016) ‘Choice of statistical model for cost-effectiveness analysis and covariate adjustment: empirical application of prominent models and assessment of their results’, *The European Journal of Health Economics* **17**:(8) 927–938.
- NICE (2013) , *Guide to the Methods of Technology Appraisal*, National Institute for Health and Care Excellence, London, UK.

- Nixon, R. M. and Thompson, S. G. (2005), ‘Methods for incorporating covariate adjustment, subgroup analysis and between-centre differences into cost-effectiveness evaluations.’ *Health economics* **14**(12), 1217–29.
- Rossi, P., Allenby, G. and McCulloch, R. (2012), *Bayesian Statistics and Marketing*, Wiley Series in Probability and Statistics, Wiley.
- Rubin, D. (1987). *Multiple Imputation for Nonresponse in Surveys*. Chichester: Wiley.
- Schmidt, P. (1990). ‘Three-stage least squares with different instruments for different equations’, *Journal of Econometrics* **43**(3), 389 – 394.
- Terza, J. V., Basu, A. and Rathouz, P. J. (2008). ‘Two-stage residual inclusion estimation: Addressing endogeneity in health econometric modeling’, *Journal of Health Economics* **27**(3), 531 – 543.
- White, I. R. and Carlin, J. B. (2010). ‘Bias and efficiency of multiple imputation compared with complete-case analysis for missing covariate values’. *Statistics in Medicine*, **28**: 2920–2931.
- White, I. R., Royston, P. and Wood, A. M. (2011). ‘Multiple imputation using chained equations: issues and guidance for practice’, *Statistics in Medicine*, **30** (4), 377–399.
- Willan, A. R. (2006). ‘Statistical Analysis of cost-effectiveness data from randomised clinical trials’. *Expert Review Pharmacoeconomics Outcomes Research* **6**, 337–346.
- Willan, A. R., Briggs, A. and Hoch, J. (2004). ‘Regression methods for covariate adjustment and subgroup analysis for non-censored cost-effectiveness data’. *Health Econ.* **13**(5), 461–475.
- Willan, A. R., Chen, E., Cook, R. and Lin, D. (2003). ‘Incremental net benefit in randomized clinical trials with qualify-adjusted survival’. *Statistics in Medicine* **22**, 353–362.
- Zellner, A. (1962). ‘An efficient method of estimating seemingly unrelated regressions and tests for aggregation bias’. *Journal of the American Statistical Association* **57**(298), pp. 348–368.
- Zellner, A. and Huang D. S. (1962). ‘Further Properties of Efficient Estimators for Seemingly Unrelated Regression Equations’. *International Economic Review* **3**(3), pp. 300–313.
- Zellner, A. and Theil, H. (1962), ‘Three-stage least squares: Simultaneous estimation of simultaneous equations’. *Econometrica* **30**(1), pp. 54–78.
- Zhang, Z., Peluso, M. J., Gross, C. P., Viscoli, C. M. and Kernan, W. N. (2014). ‘Adherence reporting in randomized controlled trials’. *Clinical Trials* **11**(2), 195–204.