



## Simulation-Based Power Calculations for Mixed Effects Modeling: **ipdpower** in **Stata**

**Evangelos Kontopantelis**

NIHR School for Primary Care Research  
University of Manchester

**David A Springate**

University of Manchester

**Rosa Parisi**

University of Manchester

**David Reeves**

University of Manchester

---

### Abstract

Simulations are a practical and reliable approach to power calculations, especially for multi-level mixed effects models where the analytic solutions can be very complex. In addition, power calculations are model-specific and multi-level mixed effects models are defined by a plethora of parameters. In other words, model variations in this context are numerous and so are the tailored algebraic calculations. This article describes **ipdpower** in **Stata**, a new simulations-based command that calculates power for mixed effects two-level data structures. Although the command was developed having individual patient data meta-analyses and primary care databases analyses in mind, where patients are nested within studies and general practices respectively, the methods apply to any two-level structure.

*Keywords:* **Stata**, **ipdpower**, power, coverage, meta analysis, multi level, mixed effects, random effects, individual patient data, IPD, primary care databases, PCD.

---

## 1. Introduction

The size of primary care databases (PCDs) allows for investigations that cannot normally be undertaken in much smaller randomized controlled trials (RCTs), such as the moderating effect of a patient characteristic on the effect of an intervention. However, researchers quite often underestimate the numbers needed to detect such effects and assume that the size of the database alone guarantees adequate power for any type of investigation. The essential structure of a PCD dataset closely resembles that of an individual patient data (IPD) meta-

analysis, with patients nested (clustered) in general practices in the former and in studies in the latter, thus the same modeling approaches and power analysis considerations are applicable to both. In individual patient data meta-analyses, the quantity of the collected data from numerous RCTs sometimes deceives researchers with regards to the power in their analysis and can lead to the assumption that a power calculation is not necessary. However, even though power calculations are quite often needed in these contexts, especially for detecting a moderating effect, available software options are not user friendly or can only cope with simple models (e.g., without random effects for the higher level clusters). Hence, it is not uncommon for researchers to use a rough ‘four-times as many patients needed as for the main effect’ over-simplification to estimate power (McClelland and Judd 1993).

More recently, an algebraic approximation approach, which uses study summary statistics, was developed by Kovalchik and Cumberland (2012) and implemented in R (Kovalchik 2013; R Core Team 2016) for individual patient data meta-analysis. Overall, Kovalchik and Cumberland found that their computationally cheap method performed well in simulations although it is constrained by distributional assumptions and performance deteriorated in some scenarios. Although this approach can be very useful, it is limited to a randomized controlled trial setting with a binary treatment, a patient-level covariate and a binary or continuous outcome. In addition, random effects for the intercept and intervention were modeled but other random effects (e.g., for the covariate) were not considered.

Using simulations to calculate power is a well-known but computationally expensive approach (Feiveson 2002), but with the availability of powerful computers, simulations can be a very useful alternative when it comes to complex study designs for which power equations are unavailable or prohibitively complex (Arnold, Hogan, Colford, Jr, and Hubbard 2011). For example, the `simsam` command offers a very flexible approach, allowing users to calculate power, coverage and other performance metrics under any probability model that can be programmed in Stata (Hooper 2013). However, the data structures need to be manually generated and the complexity of multi-level model structures, random-effects, interaction effects and non-normal data makes such an approach very difficult or even prohibitive. The `ipdpower` command for Stata (StataCorp. 2011) provides a simple but flexible framework for simulation-based power calculations via regression modeling, for a plethora of complex data structures. Various random effects at the cluster (e.g., practice, study) level can be hypothesized, the outcome can be continuous, binary or count data, covariates and moderators (i.e., interactions) can be modeled at the cluster- or patient-level, exposure can be binary (e.g., intervention) or continuous (as is often the case in observational studies), distributions for the random effects and continuous variables can be simulated to be non-normal, and ‘missingness’ mechanisms for the outcome can be implemented. Besides power, the command provides additional information that can be useful when designing a study, for example coverage (type I error) which might depart from the nominal level when the data structure is complex.

## 2. Methods

The command proceeds in two steps, within each simulation iteration. First, it generates a dataset, defining the outcome according to the specified coefficients for the intercept (constant), the exposure (intervention), the covariate and the exposure-covariate interaction (moderator effect). Second, it uses regression modeling to calculate model fit statistics, power and coverage for the dataset. Information is then aggregated across all simulated datasets. Power

indicates the percentage of iterations in which a model coefficient was found to be statistically significant and in the hypothesized direction. Coverage indicates the percentage of confidence intervals around the coefficient that include the true value and should correspond to the hypothesized  $\alpha$  level. Binomial proportion confidence intervals using the `ci` command are calculated for both power and coverage.

## 2.1. Dataset generation

For a continuous outcome, each dataset is generated using the following set of equations:

$$Y_{ij} = \beta_{0j} + \beta_{1j}\text{group}_{ij} + \beta_{2j}X_{ij} + \beta_{3j}\text{group}_{ij} * X_{ij} + \epsilon_{ij} \quad (1a)$$

with

$$\begin{aligned} \beta_{0j} &= \gamma_0 + u_{0j} \\ \beta_{1j} &= \gamma_1 + u_{1j} \\ \beta_{2j} &= \gamma_2 + u_{2j} \\ \beta_{3j} &= \gamma_3 + u_{3j} \end{aligned} \quad (2)$$

and when assuming a normal distribution for the error and random effects

$$\begin{aligned} \epsilon_{ij} &\sim N(0, \sigma_j^2) \\ u_{0j} &\sim N(0, \tau_0^2) \\ u_{1j} &\sim N(0, \tau_1^2) \\ u_{2j} &\sim N(0, \tau_2^2) \\ u_{3j} &\sim N(0, \tau_3^2) \end{aligned} \quad (3)$$

where

- $i$  the patient,
- $j$  the the cluster (e.g., study),
- $Y_{ij}$  the outcome for patient  $i$  in cluster  $j$ ,
- $\text{group}_{ij}$  exposure for patient  $i$  in cluster  $j$ ,
- $X_{ij}$  the covariate (e.g., baseline level) for patient  $i$  in cluster  $j$ ,
- $\text{group}_{ij} * X_{ij}$  the exposure-covariate interaction (moderator) for patient  $i$  in cluster  $j$ ,
- $\beta_{0j}$  the intercept for cluster  $j$ ,
- $\gamma_0$  the mean common intercept,
- $\beta_{1j}$  the exposure effect for cluster  $j$ ,
- $\gamma_1$  the mean exposure effect,
- $\beta_{2j}$  the covariate effect for cluster  $j$ ,

- $\gamma_2$  the mean covariate effect,
- $\beta_{3j}$  the interaction effect for cluster  $j$ ,
- $\gamma_3$  the mean interaction effect,
- $u_{0j}$  the random intercept for cluster  $j$ ,
- $u_{1j}$  the random exposure effect for cluster  $j$ ,
- $u_{2j}$  the random covariate effect for cluster  $j$ ,
- $u_{3j}$  the random interaction effect for cluster  $j$ ,
- $\tau_0^2$  the between-cluster variance for the intercept,
- $\tau_1^2$  the between-cluster variance for the exposure effect,
- $\tau_2^2$  the between-cluster variance for the covariate effect,
- $\tau_3^2$  the between-cluster variance for the interaction effect,
- $\epsilon_{ij}$  the error term for patient  $i$  in cluster  $j$ ,
- $\sigma_j^2$  the within-cluster variance for cluster  $j$ .

If the outcome is dichotomous Equation 1a becomes:

$$\ln\left(\frac{p_{ij}}{1-p_{ij}}\right) = \beta_{0j} + \beta_{1j}\text{group}_{ij} + \beta_{2j}X_{ij} + \beta_{3j}\text{group}_{ij} * X_{ij} \quad (1b)$$

where  $p_{ij}$  the probability that the outcome is 1 for patient  $i$  in cluster  $j$  and coefficients now correspond to log odds ( $\beta_0$ ) and log odds ratios ( $\beta_1, \beta_2$  and  $\beta_3$ ).

And for a count outcome, Equation 1a becomes:

$$\ln(E(Y_{ij})) = \beta_{0j} + \beta_{1j}\text{group}_{ij} + \beta_{2j}X_{ij} + \beta_{3j}\text{group}_{ij} * X_{ij} \quad (1c)$$

where  $E(Y_{ij})$  the expectation of  $Y_{ij}$  and the coefficients correspond to log incidence ( $\beta_0$ ) and log incidence rates ( $\beta_1, \beta_2$  and  $\beta_3$ ).

However, we also implemented two alternative skew-normal distributions for the random effects and residual error using methods described by [Ramberg, Dudewicz, Tadikamalla, and Mykytka \(1979\)](#). Besides normal distributions (skewness = 0; kurtosis = 3), **ipdpower** can model moderate-skew (skewness = 1; kurtosis = 4) and extreme-skew (skewness = 2; kurtosis = 9) distributions with the same means and variances as described in Equation 3. Although non-normal distributions for the random effects have not been found to notably affect the performance of two-stage meta-analysis methods ([Kontopantelis and Reeves 2012](#)), deviations from normality are not necessarily insignificant in this context.

**ipdpower** allows various random effects components (intercept, exposure, covariate and interaction) to be incorporated in each dataset. In meta-analysis nomenclature, the variance associated with random effects is described as heterogeneity of the exposure (treatment) effect in two-stage meta-analyses (stage 1: obtain or calculate study results; stage 2: calculate

overall effect using an appropriate method). When individual patient data are available, meta-analysts can choose one of many mixed effects regression models to combine information and model various random effects in a single stage (Kontopantelis and Reeves 2013), which is considered the best approach to meta-analysis (Mathew and Nordstrom 2010). Usually heterogeneity is quantified with  $I^2$  (Higgins and Thompson 2002) or  $H^2$  (Mittlbock and Heinzl 2006), using the within ( $\sigma^2$ ) and between-study ( $\tau_1^2$ ) variance estimates (assuming within-variance is similar across clusters). Since the error  $\epsilon$  (defined using variance  $\sigma_j^2$ ) is only meaningful in OLS regressions, these heterogeneity measures are only relevant when the outcome is continuous in IPD meta-analysis (in Stata,  $\sigma_j^2$  is fixed to  $\pi^2/3$  for logistic regressions). For this reason `ipdpower` accepts between-cluster variance inputs ( $\tau_0^2$ ,  $\tau_1^2$ ,  $\tau_2^2$  and  $\tau_3^2$ ) to define the random effects for each cluster. However, when the outcome is continuous and users wish to utilize  $I^2$  as a starting point to define model heterogeneity, solving for  $\tau^2$  returns (Higgins and Thompson 2002):

$$\tau_i^2 = \frac{\sigma^2 I_i^2}{100 - I_i^2} \quad (4)$$

Similarly,  $\tau_i^2 = (H_i^2 - 1)\sigma^2$ , where  $i \in \{0, 1, 2, 3\}$ , and the command also provides  $I^2$  and  $H^2$  as outputs for users to double-check their calculations.

## 2.2. Regression modeling

`ipdpower` provides seven different regression modeling options, consistent across the three possible outcome types, to account for the various levels of complexity in the generated datasets. The simplest analysis approach, `model(1)`, does not account for clustering, assumes no random effects and employs commands `regress`, `logit` or `poisson` for continuous, binary or count outcomes respectively. A more advanced approach (`model(2)`), declares the data as clustered using `xtset` and analyses with `xtreg`, `xtlogit` or `xtpoisson`. With this family of models only a random effects component for the intercept ( $\tau_0^2 > 0$ ) is considered and estimated. More advanced modeling options, which allow consideration of more random effects components, have also been implemented. However, these models which use `xtmixed`, `xtmelogit` and `xtmepoisson` (renamed in Stata 13), can be computationally expensive and do not always converge. These modeling approaches have been described in the IPD meta-analysis context (Higgins, Whitehead, Turner, Omar, and Thompson 2001; Turner, Omar, Yang, Goldstein, and Thompson 2000; Whitehead 2002; Kontopantelis and Reeves 2013). The simplest of them, `model(3)`, assumes a fixed common intercept, random exposure effects ( $\tau_1^2 > 0$ ) and fixed effect for the covariate. The recommended model for meta-analyses (Whitehead 2002), `model(4)`, assumes fixed study specific intercepts, random exposure effects ( $\tau_1^2 > 0$ ) and fixed study specific effects for the covariate. `model(5)` assumes random study intercepts ( $\tau_0^2 > 0$ ), random exposure effects ( $\tau_1^2 > 0$ ) and fixed study specific effect for the covariate. `model(6)` assumes random study intercepts ( $\tau_0^2 > 0$ ), random exposure effects ( $\tau_1^2 > 0$ ) and random effects for the covariate ( $\tau_2^2 > 0$ ). `model(7)` assumes random study intercepts ( $\tau_0^2 > 0$ ), random exposure effects ( $\tau_1^2 > 0$ ), random effects for the covariate ( $\tau_2^2 > 0$ ) and random effects for the interaction ( $\tau_3^2 > 0$ ). The more complicated the model the more likely analysis will fail to converge in a particular dataset, and this is often the case for the last three models – especially the last one.

Although fixed effect models are widely used in two-stage meta-analyses, even when heterogeneity is not zero (Kontopantelis, Springate, and Reeves 2013), accounting for even low levels

of between-cluster variability is a more conservative approach (Hunter and Schmidt 2000). When a fixed effects model is incorrectly assumed, both coverage and power deteriorate as true heterogeneity increases (Brockwell and Gordon 2001; Kontopantelis and Reeves 2012). Analogously, for patient data analyses, we would expect poor fit from `model(1)`, the fixed effect approach, in the presence of heterogeneity. More generally, the closer the modeling assumptions to the true data structure, the better the expected performance. However, as explained previously, the more complex models come with practical limitations so users must make the choice that best suits their needs in this trade-off between feasibility and model correctness. From this point of view, the usability of `ipdpower` is not limited to power calculations for a moderator effect, but it can be used to evaluate the overall performance of current modeling approaches in various simulated scenarios. For example, to compare `model(4)`, the preferred IPD meta-analysis approach, to the simpler `model(3)`.

### 3. The `ipdpower` command

#### 3.1. Syntax

```
ipdpower, sn(#) ssl(#) ssh(#) b0(#) b1(#) b2(#) b3(#) [minsh(#) hpoisson
  icluster outc(string) cb(#) cexp cexpd(string) errsd(#) sderrsd(#)
  derr(string) ccovd(string) bcov bcb(#) slcov tsq0(#) tsq1(#) tsq2(#)
  tsq3(#) dtp0(string) dtp1(string) dtp2(string) dtp3(string)
  covmat(name) missp(#) mar(#) mnar(#) minum(#) mipmm(#) model(#) clvl(#)
  seed(#) nskip dnorm xnodts nodisplay moreon ]
```

#### 3.2. Required

**sn(#)** Number (*integer*) of simulations to execute. At least 1000 are recommended for relatively narrow confidence intervals.

**ssl(#)** Total number (*integer*) of patients across all higher level units (clusters).

**ssh(#)** Number (*integer*) of higher level units (e.g., studies, general practices etc.).

**b0(#)** Coefficient (*real*) for the intercept (constant) of regression model. For logistic and Poisson regression log odds and log incidence rates are expected respectively. The coefficient can be zero.

**b1(#)** Coefficient (*real*) for the exposure variable (e.g., intervention: treatment vs no treatment) of the regression model. For logistic and Poisson regression log odds ratios and log incidence rate ratios are expected respectively. The coefficient can be zero.

**b2(#)** Coefficient (*real*) for the covariate variable (e.g., age) of the regression model. It can be continuous (default) or binary (`bcov` option), and patient-level (default) or study-level (`slcov` option). For logistic and Poisson regression log odds ratios and log incidence rate ratios are expected respectively. The coefficient can be zero. Note that the covariate, when continuous, is always assumed to be standardized (mean = 0 and sd = 1) since interactions are included in the models. Users need to take that into consideration when deciding on `b2`.

**b3(#)** Coefficient (*real*) for the exposure-covariate interaction variable. The command can automatically handle a binary by continuous, binary by binary or continuous by continuous interaction term. For logistic and Poisson regression log odds ratios and log incidence

rate ratios are expected respectively. The coefficient can be zero. Note that the covariate, when continuous, is always assumed to be standardized (mean = 0 and sd = 1) to allow for meaningful estimates. Users need to take that into consideration when deciding on `b3`.

### 3.3. Optional

#### *Data Structure*

`minsh(#)` Minimum number (*integer*) of patients in a higher level unit. The default is 50 since this is usually the threshold above which the effort required to obtain individual patient data for meta-analysis is justified. The command uses the numbers provided with `ssl(#)`, `ssh(#)` and `minsh(#)` to draw sizes for the higher level units from a uniform distribution. If the average size for the higher level unit is smaller than the minimum number of patients the command will return an error.

`hpoisson` Inform the command that the higher level unit sizes will not be drawn from a uniform but a Poisson distribution with mean = `ssl(#)/ssh(#)`. This approach provides cluster sizes that are much more similar in size. Cannot be used with option `minsh(#)`.

`icluster` Inform the command that the exposure is clustered at the higher level units (e.g., cluster-RCT). Can only be selected when exposure is binary and with balanced designs (i.e., cannot be used with `cexp` or `cb(#)`). Clusters assigned an odd identifier are assumed to include controls and clusters assigned even identifiers are assumed to include the intervention cases (i.e., if an odd number of higher level units is simulated with `icluster` there will be an additional cluster of controls).

`outc(string)` Type of outcome: continuous ('cont' default); dichotomous ('binr'); count ('count'). The model for a continuous outcome is  $y = b_0 + b_1 * grp + b_2 * xcovar + b_3 * xcovar * grp + u_0 + u_1 * grp + u_2 * xcovar + u_3 * xcovar * grp + errx$ , where `grp` the exposure variable, `xcovar` the covariate, `xcovar * grp` their interaction, `u0–u3` the random effects components and `errx` the residual errors. For a binary outcome the model is  $y = \text{uniform}() < \text{invlogit}(b_0 + b_1 * grp + b_2 * xcovar + b_3 * xcovar * grp + u_0 + u_1 * grp + u_2 * xcovar + u_3 * xcovar * grp)$  and for a count outcome it is  $y = \text{rpoisson}(\exp(b_0 + b_1 * grp + b_2 * xcovar + b_3 * xcovar * grp + u_0 + u_1 * grp + u_2 * xcovar + u_3 * xcovar * grp))$ . Negative binomial models were considered as an alternative for count data but they were thought to be too complex to be of much practical use in the assumption-laden context of power calculations. Note that residual errors can only be directly controlled in the OLS regression model.

`cb(#)` Probability (*real*) for patient membership to the exposure group (`grp = 1`), when exposure is binary. The default is 0.5 for a balanced design.

`cexp` Inform the command that exposure is continuous (standardized, i.e., mean = 0 and sd = 1) rather than binary (the default).

`cexpd(string)` Distribution for continuous exposure: normal ('norm' default); moderate skew ('sknorm'); extreme skew ('xsknorm'). Normal distribution for the exposure (skew = 0, kurtosis = 3) is the default. Moderate (skew = 1, kurtosis = 4) and extreme skewness (skew = 2, kurtosis = 9) are implemented using [Ramberg et al. \(1979\)](#) method. The distribution for the exposure will not affect the distribution of the outcome much unless `b1` is reasonably large. Note that the exposure, when continuous, is always assumed to be standardized (mean

= 0 and sd = 1) since interactions are included in the models. Users need to take that into consideration when deciding on **b1** and **b3**.

**errsd(#)** For continuous outcome only, standard deviation for the residual error (*real*). The default is 1. This value, combined with the model coefficients, will affect the model fit; e.g., a large value will drive down the average adjusted  $R^2$ . It also affects model heterogeneity since this is effectively the sd for the outcome within the higher level unit (e.g., within-study variability).

**sderrsd(#)** For continuous outcome only, standard deviation for the standard deviation of the residual error (*real*). In other words, it allows the residual error to vary across higher-level units, which might be a more realistic modeling strategy. The default is 0, not allowing variation which complies with modeling and heterogeneity assumptions (e.g., pooled within-variance is used for heterogeneity calculations).

**derr(string)** For continuous outcome only, distribution for errors and hence outcome: normal ('norm' default); moderate skew ('sknorm'); extreme skew ('xsknorm'). Normal distribution for the error (skew = 0, kurtosis = 3) is the default. Moderate (skew = 1, kurtosis = 4) and extreme skewness (skew = 2, kurtosis = 9) are implemented using [Ramberg et al. \(1979\)](#) method. The distribution for the errors will affect the distribution of the outcome and the larger the modeled errors the more similar the two distributions.

**ccovd(string)** Distribution for continuous covariate: normal ('norm' default); moderate skew ('sknorm'); extreme skew ('xsknorm'). Normal distribution for the covariate (skew = 0, kurtosis = 3) is the default. Moderate (skew = 1, kurtosis = 4) and extreme skewness (skew = 2, kurtosis = 9) are implemented using [Ramberg et al. \(1979\)](#) method. The distribution for the covariate will not affect the distribution of the outcome much unless **b2** is reasonably large. Note that the covariate, when continuous, is always assumed to be standardized (mean = 0 and sd = 1) since interactions are included in the models. Users need to take that into consideration when deciding on **b2** and **b3**.

**bcov** Inform the command that the covariate is binary covariate instead of continuous (the default).

**bcb(#)** For binary covariate only, probability (*real*) that **xcovar** = 1. The default is 0.5.

**slcov** Inform the command that the covariate is higher-level (e.g., study-level: recruitment setting) rather than patient-level(the default).

### *Random effects*

**tsq0(#)** Random effect between higher level variance for the intercept. The default value is 0, which assumes homogeneity and no random effects. Heterogeneity for this model factor (intercept) is calculated using **tsq0** and **errsd**. For example  $I^2 = 100 * \text{tsq0} / (\text{tsq0} + \text{errsd}^2)$  and  $H^2 = 1 / (1 - \text{tsq0} / (\text{tsq0} + \text{errsd}^2))$ . Solving for **tsq0** we obtain  $\text{tsq0} = (I^2 * \text{errsd}^2) / (100 - I^2)$  and  $\text{tsq0} = H^2 * \text{errsd}^2 - \text{errsd}^2$ . Although *ipdpower* does not allow  $I^2$  or  $H^2$  inputs for the random effects components, they can be easily calculated using these formulas. Additionally users can use the obtained the hypothesized heterogeneity levels for the inputted within- and between- variance parameters (say in a small trial simulation, if unsure of calculations).

**tsq1(#)** Random effect between higher level variance for the exposure. The default value is 0, which assumes homogeneity and no random effects. Heterogeneity for this model factor (exposure) would be calculated using **tsq1** and **errsd**. See **tsq0(#)** for calculation details.



`tsq2(#)` Random effect between higher level variance for the covariate. The default value is 0, which assumes homogeneity and no random effects. Heterogeneity for this model factor (covariate) would be calculated using `tsq2` and `errsd`. See `tsq0(#)` for calculation details.

`tsq3(#)` Random effect between higher level variance for the exposure \* covariate interaction term. The default value is 0, which assumes homogeneity and no random effects. Heterogeneity for this model factor (interaction) would be calculated using `tsq3` and `errsd`. See `tsq0(#)` for calculation details.

`dtp0(string)` Distribution for intercept random effect: normal (`'norm'` default); moderate skew (`'sknorm'`); extreme skew (`'xsknorm'`). Normal distribution (skew = 0, kurtosis = 3) is the default. Moderate (skew = 1, kurtosis = 4) and extreme skewness (skew = 2, kurtosis = 9) are implemented using the [Ramberg \*et al.\* \(1979\)](#) method. In the standard two-stage meta-analysis setting, a non-normal distribution for the random effects has been found to have a small effect on power and coverage – even when the distribution is quite extreme.

`dtp1(string)` Distribution for exposure random effect: normal (`'norm'` default, skew = 0 kurtosis = 3); moderate skew (`'sknorm'`, skew = 1 kurtosis = 4); extreme skew (`'xsknorm'`, skew = 2 kurtosis = 9).

`dtp2(string)` Distribution for covariate random effect: normal (`'norm'` default, skew = 0 kurtosis = 3); moderate skew (`'sknorm'`, skew = 1 kurtosis = 4); extreme skew (`'xsknorm'`, skew = 2 kurtosis = 9).

`dtp3(string)` Distribution for exposure \* covariate interaction random effect: normal (`'norm'` default, skew = 0 kurtosis = 3); moderate skew (`'sknorm'`, skew = 1 kurtosis = 4); extreme skew (`'xsknorm'`, skew = 2 kurtosis = 9).

`covmat(name)` Covariance matrix for normally distributed random effects. Alternative random effects definition to allow modeling of relationships between the random effects components. The matrix needs to be a 4x4 symmetrical non-negative matrix, with the diagonal elements corresponding to the random effects variances for intercept (`name [1,1]`), exposure (`name [2,2]`), covariate (`name [3,3]`) and interaction (`name [4,4]`). Non-normal random effects cannot be modeled using this approach.

### *Missing data*

`missp(#)` Probability (*real*) that outcome is missing, to allow for missing data mechanisms. If option `mar(#)` is defined, data are assumed to be missing under a missing at random (MAR) mechanism. If option `mnar(#)` is defined, data are assumed to be missing under a missing not at random (MNAR) mechanism. If neither `mar(#)` or `mnar(#)` are provided along with `missp(#)`, data are assumed to be missing under a missing completely at random (MCAR) mechanism. Please note that multiple imputation models can satisfactorily deal with MCAR and MAR mechanisms and not MNAR mechanisms, although they will improve power in all three scenarios. In terms of minimizing estimate bias, multiple imputations are not really needed for MCAR data (a complete case analysis should provide similar estimates) and are mainly needed for MAR data. There is some evidence that multiple imputation can offer some protection against MNAR mechanisms, but in general the estimates from such models will be biased. However, it is impossible to assess whether data are MNAR, without obtaining additional external data, and the multiple imputation models are now used routinely whenever ‘missingness’ is encountered. Therefore, we decided to offer a MNAR modeling option with the under-performing, in this scenario, multiple imputation approach to allow inquisitive

researchers to practically calculate the performance of these models, in terms of power.

**mar(#)** Odds ratio (*real*) that defines a missing at random (MAR) mechanism. The relationship between the covariate and missingness in the outcome is defined by  $z = \ln(\text{mar}(\#)) * \text{xcovar}$  (i.e., a logistic regression model), and the missing data are selected (i.e., set to missing) from the  $z = 1$  sub-sample. So a value of 1 implies the mechanism is MCAR, a value above 1 implies that the outcome is more likely to be missing for larger values of the covariate and a value below 1 implies that the outcome is more likely to be missing for smaller values of the covariate.

**mnar(#)** Odds ratio (*real*) that defines a missing not at random (MNAR) mechanism. The relationship between the outcome and missingness in the outcome is defined by  $z = \ln(\text{mnar}(\#)) * y$  (i.e., a logistic regression model), and the missing data are selected (i.e., set to missing) from the  $z = 1$  sub-sample. So a value of 1 implies the mechanism is MCAR, a value above 1 implies that the outcome is more likely to be missing for larger values of the outcome and a value below 1 implies that the outcome is more likely to be missing for smaller values of the outcome.

**minum(#)** Number (*integer* > 1) of multiple imputations to be executed. This options informs *ipdpower* that multiple-imputation models will be used and therefore missing data with one of the three available structures (MCAR, MAR, MNAR) need to have been defined. For the imputations, univariate linear, logistic or Poisson regression is used depending on the outcome (see *mi impute*). Under all imputation models, the outcome depends on exposure, covariate and their interaction, and for multi-level models (2–7) additionally on higher level units (clusters). The imputed datasets are then analyzed using *mi estimate* as a prefix, for the seven available models. Note that this process can be time consuming for complex models and binary or count outcomes, while convergence issues are amplified since all imputed datasets must run successfully for *mi estimate* to return results. Therefore 5 imputations are recommended for most models, and 2–3 for non-continuous outcomes and models 5, 6 or 7 (see *model(#)* below).

**mipmm(#)** For a continuous outcome, it informs that missing data will be imputed using a predictive mean matching algorithm (rather than linear regression). The algorithm is computationally more expensive and # defines the number (*integer* ≥ 1) of closest observations (nearest neighbors) to draw from.

### *Modeling*

**model(#)** Specifies the form of the regression model: 1 simple, 2 random effects for intercept, 3–7 various mixed effects options. Model 1 corresponds to a regression with *regress*, *logit* or *poisson*, for continuous, binary and count outcomes respectively. Random effects are not considered at all under these models. Model 2 uses the *xt* family of models, sets the higher level as a panel variable with *xtset* and analyses with *xtreg*, *xtlogit* or *xtpoisson*. Only a random effects component for the intercept is considered under with this set of models. Models 3–7 allow for more advanced modeling options, accounting for various random effect components, but are computationally expensive and do not always converge (using *xtmixed*, *xtmelogit* and *xtmepoisson* which have been renamed in Stata 13 – but we wished to ensure *ipdpower* was compatible with Stata 12). The modeling approaches have been described for *ipdforest* and in the following descriptions we assume the highel level is study (i.e., patients nested within studies), for convenience. Model 3 assumes a fixed common intercept, random

#	Commands	Intercept	Exposure	Covariate	Interaction
1	<code>regress/logit/poisson</code>	Fixed	Fixed	Fixed	Fixed
2	<code>xtreg/xtlogit/xtpoisson</code>	Random	Fixed	Fixed	Fixed
3	<code>xtmixed/xtmelogit/xtmepoisson</code>	Fixed	Random	Fixed	Fixed
4	<code>xtmixed/xtmelogit/xtmepoisson</code>	Fixed <sup>†</sup>	Random	Fixed <sup>†</sup>	Fixed
5	<code>xtmixed/xtmelogit/xtmepoisson</code>	Random	Random	Fixed <sup>†</sup>	Fixed
6	<code>xtmixed/xtmelogit/xtmepoisson</code>	Random	Random	Random	Fixed
7	<code>xtmixed/xtmelogit/xtmepoisson</code>	Random	Random	Random	Random

Table 1: Modeling options with **ipdpower** (<sup>†</sup>cluster-specific fixed-effects; different estimate for each higher level unit rather than overall)

exposure effects and fixed effect for the covariate. Model 4 assumes fixed study specific intercepts, random exposure effects and fixed study specific effects for the covariate (which is usually the recommended model for performing individual patient data meta-analysis). Model 5 assumes random study intercepts, random exposure effects and fixed study specific effects for the covariate. Model 6 assumes random study intercepts, random exposure effects and random effects for the covariate. Model 7 assumes random study intercept, random exposure effects, random effects for the covariate and random effects for the interaction. Models 5 and 6 often fail to converge and for model 7 non-convergence is more frequent than convergence. The models are summarized in Table 1.

`clvl(#)` Set confidence level. The default is 95% (alpha level of 5%). See `level`.

`seed(#)` Set initial value of random-number seed, for the simulations. See `set seed`.

`nskip` Add `noskip` option to `xtlogit` or `xtpoisson` to return McFadden's pseudo  $R^2$ . This option is only relevant when `model(2)` with `outc(1)` or `outc(2)` is used. Computationally, this approach is more expensive since it additionally fits a full maximum-likelihood model with only a constant for the regression equation be fit (which is used as the base model for the comparison with the final model).

`dnorm` Add `normal` option to `xtpoisson` to assume normally distributed random effects for the intercept. The default is gamma-distributed which is computationally less expensive. Additionally, when the `normal` option is specified, model convergence often fails. This option is only relevant when `model(2)` with `outc(2)`. A skew-normal distribution is similar to gamma and perhaps skew-normal random effects should be considered when modeling count data.

`xnodts` Suppress simulation progress display. If option not specified, a '.' is displayed for each successful model run (i.e., converging) and an 'x' for each unsuccessful iteration.

`nodisplay` Do not display results at the end of the simulation process. Suppressed results include: simulation characteristics, average model fit, average statistics for the outcome, average `b0–b3`, hypothesized heterogeneity values, power and coverage.

`moreon` Set `more on` (default is off).

### 3.4. Saved results

Table 2 lists the saved scalar results of `ipdpower` in `r()`. Coefficient estimates, power, coverage, simulations and computational time information is always returned. An  $R^2$  statistic

Name	Description
r(bo)	Average coefficient estimate for the intercept
r(b1)	Average coefficient estimate for the exposure
r(b2)	Average coefficient estimate for the covariate
r(b3)	Average coefficient estimate for the interaction
r(nsim)	Number of simulations
r(nrun)	Number of successful simulations
r(ctime)	Computational time (minutes)
r(rsq)	Average adjusted or pseudo $R^2$
r(errsd)	Within-sd (error)
r(consd)	between-sd (intercept)
r(grpsd)	Within-sd (exposure)
r(covsd)	between-sd (covariate)
r(intsd)	Within-sd (interaction)
r(pow0)	Power to detect b0
r(lpow0)	Power to detect b0, lower CI
r(upow0)	Power to detect b0, upper CI
r(pow1)	Power to detect b1
r(lpow1)	Power to detect b1, lower CI
r(upow1)	Power to detect b1, upper CI
r(pow2)	Power to detect b2
r(lpow2)	Power to detect b2, lower CI
r(upow2)	Power to detect b2, upper CI
r(pow3)	Power to detect b3
r(lpow3)	Power to detect b3, lower CI
r(upow3)	Power to detect b3, upper CI
r(cov0)	Coverage for b0
r(lcov0)	Coverage for b0, lower CI
r(ucov0)	Coverage for b0, upper CI
r(cov1)	Coverage for b1
r(lcov1)	Coverage for b1, lower CI
r(ucov1)	Coverage for b1, upper CI
r(cov2)	Coverage for b2
r(lcov2)	Coverage for b2, lower CI
r(ucov2)	Coverage for b2, upper CI
r(cov3)	Coverage for b3
r(lcov3)	Coverage for b3, lower CI
r(ucov3)	Coverage for b3, upper CI

Table 2: Saved results of `ipdpower` in `r()` (CI = confidence interval).

cannot be returned for the more advanced multi-level models with `xtmixed`, `xtmelogit` or `xtnepoisson` or for multiple imputations. With `model(1)`, the average adjusted  $R^2$  is returned from `regress` and the average pseudo  $R^2$  from `logit` or `poisson`. With `model(2)`, the average overall  $R^2$  is returned from `xtreg` and, if the `nskip` option is specified, the average overall McFadden's pseudo  $R^2$  from `xtlogit` or `xtpoisson`. Within-sd is only reported for

continuous outcomes and when data are analyzed with `xtreg` or `xtmixed`. The between-sd components are only returned if accounted for and estimated by the specified model.

### 3.5. Example

As an example, we used `ipdpower` to calculate the power in a few designs where the outcome is continuous. First we assumed 20 clusters, a total of 5000 patients, no random effects and the default level for the residual error  $\epsilon$  ( $\sigma = 1$ ). The generated outcome for the model could be described by  $Y_{ij} = 1 + 0.5\text{group}_{ij} + 0.3X_{ij} + 0.1\text{group}_{ij} * X_{ij} + \epsilon_{ij}$  and we proceeded to analyze using option `model(2)`, i.e., `xtreg` allowing for a random intercept.

```
. ipdpower, sn(1000) ssl(5000) ssh(20) b0(1) b1(0.5) b2(0.3) b3(0.1) model(2)
> xnodts seed(7)
```

```
model 2: random effects regression, for study i.e. intercept
outcome type:  continuous
exposure type:  binary
covariate type: continuous
random seed number:          7
number of converging runs: 1000
computational time (min):    1.8
```

Characteristics for the outcome

```
-----
              |  group0   group1
-----+-----
mean          |    0.983   1.512
sd            |    1.038   1.104
-----
```

Modelled variance and heterogeneity measures

```
-----
              |  exposure  covariate  interaction  intercept
-----+-----
between variance (tau^2) |    0.000    0.000    0.000    0.000
I^2 (range: 0 to 100%)  |    0.000    0.000    0.000    0.000
H^2 (range: 1 to +inf)  |    1.000    1.000    1.000    1.000
-----
```

modelled within-study variance (pooled): 1.000

Results: model estimates

```
-----
              |  exposure  covariate  interaction  intercept
-----+-----
coefficient mean |    0.500    0.300    0.101    1.000
between-sd       |    .         .         .         0.013
-----
```

```
within-sd(error):    1.000
R^2(%):              15.834
```

Results: coverage

	estimate	[95% Conf. Interval]	
exposure	94.2	92.6	95.6
covariate	94.6	93.0	95.9
interaction	95.9	94.5	97.0
intercept	96.1	94.7	97.2

Results: power

	estimate	[95% Conf. Interval]	
exposure	100.0	99.6	100.0
covariate	100.0	99.6	100.0
interaction	95.4	93.9	96.6
intercept	100.0	99.6	100.0

The model setup details (within/between-variance, heterogeneity, outcome characteristics etc.) are always provided as outputs for confirmation. On average, the model explained approximately 15.8% of the variance and performed well in the parameter estimation. Assuming we are interested in capturing the hypothesized interaction effect, the power to detect it was estimated at 95.4%. Coverage for the interaction effect was above nominal at 95.9%.

In the next step we introduce heterogeneity for the exposure.

```
. ipdpower, sn(1000) ssl(5000) ssh(20) b0(1) b1(0.5) b2(0.3) b3(0.1)
> tsq1(0.5) model(2) xnodts seed(7)
```

model 2: random effects regression, for study i.e. intercept

```
outcome type:    continuous
exposure type:   binary
covariate type:  continuous
random seed number:      7
number of converging runs: 1000
computational time (min): 3.8
```

Characteristics for the outcome

	group0	group1
mean	1.026	1.513
sd	1.045	1.278

---

 Modelled variance and heterogeneity measures
 

---

	exposure	covariate	interaction	intercept
between variance ( $\tau^2$ )	0.500	0.000	0.000	0.000
$I^2$ (range: 0 to 100%)	33.333	0.000	0.000	0.000
$H^2$ (range: 1 to +inf)	1.500	1.000	1.000	1.000

---

 modelled within-study variance (pooled): 1.000

## Results: model estimates

	exposure	covariate	interaction	intercept
coefficient mean	0.504	0.300	0.100	0.998
between-sd	.	.	.	0.347

within-sd(error): 1.057

 $R^2(\%)$ : 13.697

## Results: coverage

	estimate	[95% Conf. Interval]	
exposure	24.7	22.1	27.5
covariate	94.9	93.3	96.2
interaction	94.9	93.3	96.2
intercept	100.0	99.6	100.0

## Results: power

	estimate	[95% Conf. Interval]	
exposure	99.6	99.0	99.9
covariate	100.0	99.6	100.0
interaction	92.2	90.4	93.8
intercept	100.0	99.6	100.0

The modeled heterogeneity was moderate, with  $I^2 = 33.3\%$ . On average, the model explained less of the variance in the outcome, approximately 13.7%. The model betas and the residual error (within-sd) were estimated accurately on average but the model assumes a random intercept and underestimates the hypothesized heterogeneity (since the selected modeling approach assumes a random intercept, when a random exposure effect exists in the simulated

data). The power to detect the interaction was slightly lower, estimated at 92.2%, but coverage was almost at the nominal level with 94.9%. However, notice the very low coverage for the exposure which was due to the introduced heterogeneity for that factor.

Next, we change the distribution for the errors from normal to extreme skew-normal.

```
. ipdpower, sn(1000) ssl(5000) ssh(20) b0(1) b1(0.5) b2(0.3) b3(0.1)
> tsq1(0.5) model(2) derr(xsknorm) xnodts seed(7)
```

```
model 2: random effects regression, for study i.e. intercept
outcome type: continuous
exposure type: binary
covariate type: continuous
random seed number: 7
number of converging runs: 1000
computational time (min): 5.8
```

Characteristics for the outcome

	group0	group1
mean	1.019	1.587
sd	1.052	1.269

Modelled variance and heterogeneity measures

	exposure	covariate	interaction	intercept
between variance ( $\tau^2$ )	0.500	0.000	0.000	0.000
$I^2$ (range: 0 to 100%)	33.333	0.000	0.000	0.000
$H^2$ (range: 1 to +inf)	1.500	1.000	1.000	1.000

modelled within-study variance (pooled): 1.000

Results: model estimates

	exposure	covariate	interaction	intercept
coefficient mean	0.502	0.300	0.098	1.000
between-sd	.	.	.	0.349

within-sd(error): 1.058

$R^2$ (%): 13.603



Results: coverage

	estimate	[95% Conf. Interval]	
exposure	25.6	22.9	28.4
covariate	95.4	93.9	96.6
interaction	94.2	92.6	95.6
intercept	100.0	99.6	100.0

Results: power

	estimate	[95% Conf. Interval]	
exposure	99.4	98.7	99.8
covariate	100.0	99.6	100.0
interaction	90.3	88.3	92.1
intercept	100.0	99.6	100.0

Although the non-normal errors did not appear to affect the model fit, power and coverage were slightly affected. The power to detect the interaction dropped to 90.3% and coverage was 94.2%. Coverage was still very problematic for the exposure.

Next, we proceeded to analyze with a more appropriate model (`model(3)`), one that accounts for the exposure heterogeneity.

```
. ipdpower, sn(1000) ssl(5000) ssh(20) b0(1) b1(0.5) b2(0.3) b3(0.1)
> tsq1(0.5) model(3) derr(xsknorm) xnodts seed(7)
```

```
model 3: fixed common intercept; random treatment effect; fixed effect for
> baseline
```

```
outcome type: continuous
exposure type: binary
covariate type: continuous
random seed number: 7
number of converging runs: 1000
computational time (min): 12.7
```

Characteristics for the outcome

	group0	group1
mean	1.019	1.587
sd	1.052	1.269

## Modelled variance and heterogeneity measures

	exposure	covariate	interaction	intercept
between variance (tau <sup>2</sup> )	0.500	0.000	0.000	0.000
I <sup>2</sup> (range: 0 to 100%)	33.333	0.000	0.000	0.000
H <sup>2</sup> (range: 1 to +inf)	1.500	1.000	1.000	1.000

modelled within-study variance (pooled): 1.000

## Results: model estimates

	exposure	covariate	interaction	intercept
coefficient mean	0.503	0.301	0.099	1.000
between-sd	0.685	.	.	.

within-sd(error): 1.000  
R<sup>2</sup>(%): .

## Results: coverage

	estimate	[95% Conf. Interval]
exposure	94.2	92.6 95.6
covariate	95.5	94.0 96.7
interaction	94.0	92.3 95.4
intercept	95.0	93.5 96.3

## Results: power

	estimate	[95% Conf. Interval]
exposure	87.1	84.9 89.1
covariate	100.0	99.6 100.0
interaction	93.2	91.5 94.7
intercept	100.0	99.6 100.0

The model performed well and slightly under-estimated the hypothesized between-study variance ( $\sqrt{0.5} \approx 0.707$ ). The power to detect the interaction was 93.2% and coverage was 94.0%. As expected, coverage for the exposure was not an issue for this model, since it is a much more accurate reflection of the hypothesized data structure.

Finally, we analyzed the data with the recommended model for IPD meta-analyses which is computationally more expensive (`model(4)`). Under this model, fixed cluster specific effects for the covariate and fixed cluster specific intercepts were assumed (i.e., different parameter

estimates for each cluster) and therefore covariate and intercept information could not be reported.

```
. ipdpower, sn(1000) ssl(5000) ssh(20) b0(1) b1(0.5) b2(0.3) b3(0.1)
> tsq1(0.5) model(4) derr(xsknorm) xnodts seed(7)
```

```
model 4: fixed study specific intercepts; random treatment effect; fixed
> study specific effect for baseline
```

```
outcome type: continuous
exposure type: binary
covariate type: continuous
random seed number: 7
number of converging runs: 1000
computational time (min): 44.4
```

Characteristics for the outcome

```
-----
          |  group0   group1
-----+-----
mean      |    1.019   1.587
sd        |    1.052   1.269
-----
```

Modelled variance and heterogeneity measures

```
-----
          |  exposure  covariate  interaction  intercept
-----+-----
between variance (tau^2) |    0.500    0.000    0.000    0.000
I^2 (range: 0 to 100%)  |   33.333    0.000    0.000    0.000
H^2 (range: 1 to +inf)  |    1.500    1.000    1.000    1.000
-----
```

modelled within-study variance (pooled): 1.000

Results: model estimates

```
-----
          |  exposure  covariate  interaction  intercept
-----+-----
coefficient mean      |    0.503      .    0.099      .
between-sd            |    0.677      .      .      .
-----
```

within-sd(error): 0.996

R^2(%): .

Results: coverage

	estimate	[95% Conf. Interval]	
exposure	93.7	92.0	95.1
covariate	.	.	.
interaction	93.6	91.9	95.0
intercept	.	.	.

Results: power

	estimate	[95% Conf. Interval]	
exposure	87.9	85.7	89.9
covariate	.	.	.
interaction	93.2	91.5	94.7
intercept	.	.	.

The model did not appear to perform better than `model(3)`. The average heterogeneity estimate was similar and so were power and coverage for the interaction at 93.2% and 93.6% respectively. However, a more thorough investigation would be required to make a convincing recommendation on model choice.

## 4. Discussion

We aimed to describe `ipdpower`, a new simulation-based power calculation command. The command can be viewed as a tool that uses simulation to evaluate a particular model under working assumption but, mainly, calculate power once a plausible model has been chosen. It offers a plethora of modeling choices for researchers, and can help them decide whether to obtain or extract data from existing sources and perform an analysis. Its novelty lies with non-normal distribution options and the various random effects assumptions that can be implemented. However, `ipdpower` is flexible and, at its simplest, the higher-level inputs can be ignored to produce power calculation of a one-level model. In addition, if a parameter is not needed (say the covariate or the interaction) the user can just set the respective coefficient to zero when defining the model structure.

Usually researchers wish to input the desired power level and estimate the number of patients and second level units required. Such an approach would be more complicated to implement in `ipdpower` since the parameters of interest are two (number of patients and number of higher-level units) and the combinations that would provide the desired power are many, unless one of them is fixed. Nevertheless, users with a little programming experience should call the command from a binary search algorithm (since power is a monotonic function of sample size), that dichotomizes the search area to arrive to the solution quickly, in order to calculate the sample sizes needed to have a desired power level (Cormen, Leiserson, Rivest, Stein, and others 2001). Either the average number of patients in a cluster or the number of clusters

should be fixed in this approach. Similarly, researchers interested in generating datasets with the specified characteristics, rather than calculating power, can easily obtain them by taking advantage of the fact that `ipdpower` keeps the last simulated dataset in memory. Calling `ipdpower` from within a loop, with `sim(1)` and a different seed number for each iteration, and appending the generated datasets (possibly adding an identifier for each iteration) will produce a sample of datasets.

It might appear confusing that the command accepts variance inputs for the between-cluster heterogeneity and returns standard deviation estimates. In the meta-analysis literature, heterogeneity is usually defined using the variance, but, on the other hand, `Stata` models typically return standard deviations. We attempted to satisfy both practices.

When modeling continuous predictors (e.g., covariate or exposure), `ipdpower` generates them as standardized (mean = 0, sd = 1) to avoid potential complications due to the interaction term. However, in most cases users will wish to hypothesize effects for non-standardized variables plus it might not be straightforward to define mean levels for outcomes through `b0(#)` (e.g.,  $\Pr(Y = 1 | \text{group} = 0)$  for a binary outcome). To help researchers in setting up the correct design we have provided an accompanying Microsoft Excel file with examples across all possible exposure-covariate scenarios, with or without interaction effects (from [http://www.statanalysis.co.uk/files/defining\\_betas.xlsx](http://www.statanalysis.co.uk/files/defining_betas.xlsx)). The command also returns detailed information for inputs, assumed heterogeneity levels and aggregated model estimates (coefficients, coverage and power). Since advanced usage of the command can be challenging, users should utilize these outputs to ensure they are modeling the desired structure, and re-specify if necessary. Regarding setting heterogeneity levels and their interpretation, there are numerous practical guides that can prove helpful (Fletcher 2007; Higgins, Thompson, Deeks, and Altman 2003). Although we present information on data missingness mechanisms in the help file, interested users can find more details in an excellent overview provided by Horton and Kleinman (2007).

When many random effects are modeled, assuming that all correlations between them are zero might not always be a realistic assumption. From a practical point of view, these correlations might have very small effects on performance while their hypothesized values for a particular study could be anyone's guess. Nevertheless, we have allowed relationships between normally distributed effects to be specified using a covariance matrix (option `covmat(name)`), and modeled with the `drawnorm` command. Such an approach was not possible when modeling non-normal random effects, however, and all correlations between non-normal effects are assumed to be zero.

`ipdpower` uses some of the mixed effects modeling commands that were renamed in `Stata 13`; `xtmixed` to `mixed`, `xtmelogit` to `meqrlogit` and `xtmepoisson` to `meqrpoisson`. However, we wished to ensure `Stata 12` users would have access to the command and our choice does not affect users of later versions. In terms of model choice, users should carefully consider the implications of using very complex models where non-convergence is not that uncommon, especially for analyses of a small number of patients and a few clusters. The command output provides the number of converging runs which can inform on modeling decisions. For example, if the rate of non-converging iterations is high, users should consider that they may be unable to use the specified modeling approach with their data, when they have collected them. In this case, they should probably consider an alternative, simpler modeling approach. On the other hand, if the selected modeling approach is considered essential (for example, to model and investigate specific random-effects), users should consider increasing the sample sizes to

improve the convergence rates.

Future work could involve using `ipdpower` as a platform to inform on the effectiveness of multiple-imputation mechanisms, especially in Primary Care Databases under MNAR assumptions. In addition, although one-stage IPD meta-analysis is considered to be more robust from a theoretical perspective (Mathew and Nordstrom 2010), the simplicity of a two-stage IPD analysis and circumstantial evidence where results were very close with the two approaches make the latter appealing to researchers. However, the asymptotic nature of the methods implies that a theoretical comparison does not provide a complete picture and a simulation study is also needed to compare the two approaches, especially for small numbers of meta-analyses studies and in the presence of heterogeneity. Until such comparisons are made we must warn researchers that heterogeneity estimates in two-stage aggregate meta-analyses can be very inaccurate and more often than not seem to fail to account for existing heterogeneity (Kontopantelis and Reeves 2012; Kontopantelis *et al.* 2013), and these methodological problems are likely to be an issue in two-stage IPD meta-analyses as well.

## Acknowledgments

This study was funded by the National Institute for Health Research (NIHR) School for Primary Care Research (SPCR), under the title ‘An analytical framework for increasing the efficiency and validity of research using primary care databases’ (Project no. 211). This paper presents independent research funded by the National Institute for Health Research (NIHR). The views expressed are those of the authors and not necessarily those of the NHS, the National Institute for Health Research or the Department of Health.

We would like to thank the two anonymous reviewers whose comments improved the manuscript and the command significantly.

## Contribution

EK designed and developed the command and wrote the manuscript. DS, RP and DR critically commented on both the manuscript and the functionality of the command.

## References

- Arnold BF, Hogan DR, Colford, Jr JM, Hubbard AE (2011). “Simulation Methods to Estimate Design Power: An Overview for Applied Research.” *BMC Medical Research Methodology*, **11**, 94. doi:10.1186/1471-2288-11-94.
- Brockwell SE, Gordon IR (2001). “A Comparison of Statistical Methods for Meta-Analysis.” *Statistics in Medicine*, **20**(6), 825–840. doi:10.1002/sim.650.
- Cormen TH, Leiserson CE, Rivest RL, Stein C, others (2001). *Introduction to Algorithms*, volume 2. MIT Press, Cambridge.
- Feiveson AH (2002). “Power by Simulation.” *Stata Journal*, **2**(2), 107–124.

- Fletcher J (2007). “Clinical Epidemiology Notes: What Is Heterogeneity and Is It Important?” *BMJ: British Medical Journal*, **334**(7584), 94. doi:10.1136/bmj.39057.406644.68.
- Higgins J, Thompson SG (2002). “Quantifying Heterogeneity in a Meta-Analysis.” *Statistics in Medicine*, **21**(11), 1539–1558. doi:10.1002/sim.1186.
- Higgins JPT, Thompson SG, Deeks JJ, Altman DG (2003). “Measuring Inconsistency in Meta-Analyses.” *BMJ: British Medical Journal*, **327**(7414), 557. doi:10.1136/bmj.327.7414.557.
- Higgins JPT, Whitehead A, Turner RM, Omar RZ, Thompson SG (2001). “Meta-Analysis of Continuous Outcome Data from Individual Patients.” *Statistics in Medicine*, **20**(15), 2219–2241. doi:10.1002/sim.918.
- Hooper R (2013). “Versatile Sample-Size Calculation Using Simulation.” *Stata Journal*, **13**(1), 21–38.
- Horton NJ, Kleinman KP (2007). “Much Ado about Nothing: A Comparison of Missing Data Methods and Software to Fit Incomplete Data Regression Models.” *The American Statistician*, **61**(1). doi:10.1198/000313007x172556.
- Hunter JE, Schmidt FL (2000). “Fixed Effects vs. Random Effects Meta-Analysis Models: Implications for Cumulative Research Knowledge.” *International Journal of Selection and Assessment*, **8**(4), 275–292. doi:10.1111/1468-2389.00156.
- Kontopantelis E, Reeves D (2012). “Performance of Statistical Methods for Meta-Analysis When True Study Effects Are Non-Normally Distributed: A Simulation Study.” *Statistical Methods in Medical Research*, **21**(4), 409–426. doi:10.1177/0962280210392008.
- Kontopantelis E, Reeves D (2013). “A Short Guide and a Forest Plot Command (`ipdforest`) for One-Stage Meta-Analysis.” *Stata Journal*, **13**(3), 574–587.
- Kontopantelis E, Springate DA, Reeves D (2013). “A Re-Analysis of the Cochrane Library Data: The Dangers of Unobserved Heterogeneity in Meta-Analyses.” *PloS One*, **8**(7), e69930. doi:10.1371/journal.pone.0069930.
- Kovalchik SA (2013). “Aggregate-Data Estimation of an Individual Patient Data Linear Random Effects Meta-Analysis with a Patient Covariate-Treatment Interaction Term.” *Biostatistics*, **14**(2), 273–283. doi:10.1093/biostatistics/kxs035.
- Kovalchik SA, Cumberland WG (2012). “Using Aggregate Data to Estimate the Standard Error of a Treatment-Covariate Interaction in an Individual Patient Data Meta-Analysis.” *Biometrical Journal*, **54**(3), 370–384. doi:10.1002/bimj.201100167.
- Mathew T, Nordstrom K (2010). “Comparison of One-Step and Two-Step Meta-Analysis Models Using Individual Patient Data.” *Biometrical Journal*, **52**(2), 271–287. doi:10.1002/bimj.200900143.
- McClelland GH, Judd CM (1993). “Statistical Difficulties of Detecting Interactions and Moderator Effects.” *Psychological Bulletin*, **114**(2), 376–390. doi:10.1037/0033-2909.114.2.376.

Mittlbock M, Heinzl H (2006). “A Simulation Study Comparing Properties of Heterogeneity Measures in Meta-Analyses.” *Statistics in Medicine*, **25**(24), 4321–4333. doi:10.1002/sim.2692.

Ramberg JS, Dudewicz EJ, Tadikamalla PR, Mykytka EF (1979). “A Probability Distribution and Its Uses in Fitting Data.” *Technometrics*, **21**(2), 201–214. doi:10.1080/00401706.1979.10489750.

R Core Team (2016). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

StataCorp (2011). *Stata Statistical Software: Release 12*. StataCorp LP, College Station. URL <http://www.stata.com/>.

Turner RM, Omar RZ, Yang M, Goldstein H, Thompson SG (2000). “A Multilevel Model Framework for Meta-Analysis of Clinical Trials with Binary Outcomes.” *Statistics in Medicine*, **19**(24), 3417–3432. doi:10.1002/1097-0258(20001230)19:24<3417::aid-sim614>3.3.co;2-c.

Whitehead A (2002). *Meta-Analysis of Controlled Clinical Trials*. John Wiley & Sons.

### **Affiliation:**

Evangelos Kontopantelis  
Centre for Health Informatics & Centre for Primary Care  
NIHR School for Primary Care Research  
*and*  
Institute of Population Health  
University of Manchester  
E-mail: [e.kontopantelis@manchester.ac.uk](mailto:e.kontopantelis@manchester.ac.uk)  
URL: <http://www.population-health.manchester.ac.uk/staff/EvanKontopantelis>

David A Springate, David Reeves  
Institute of Population Health  
University of Manchester  
E-mail: [david.springate@manchester.ac.uk](mailto:david.springate@manchester.ac.uk), [david.reeves@manchester.ac.uk](mailto:david.reeves@manchester.ac.uk)  
URL: <http://www.population-health.manchester.ac.uk/staff/daspringate>  
<http://www.population-health.manchester.ac.uk/staff/76561>



Rosa Parisi  
Manchester Pharmacy School  
University of Manchester  
E-mail: [rosa.palisi@manchester.ac.uk](mailto:rosa.palisi@manchester.ac.uk)  
URL: <http://www.pharmacy.manchester.ac.uk/staff/rosapalisi>