

SCIENTIFIC REPORTS



OPEN

Using instrumental variables to disentangle treatment and placebo effects in blinded and unblinded randomized clinical trials influenced by unmeasured confounders

Elias Chaibub Neto

Received: 07 July 2016

Accepted: 25 October 2016

Published: 21 November 2016

Clinical trials traditionally employ blinding as a design mechanism to reduce the influence of placebo effects. In practice, however, it can be difficult or impossible to blind study participants and unblinded trials are common in medical research. Here we show how instrumental variables can be used to quantify and disentangle treatment and placebo effects in randomized clinical trials comparing control and active treatments in the presence of confounders. The key idea is to use randomization to separately manipulate treatment assignment and psychological encouragement conversations/interactions that increase the participants' desire for improved symptoms. The proposed approach is able to improve the estimation of treatment effects in blinded studies and, most importantly, opens the doors to account for placebo effects in unblinded trials.

Placebo effects have drawn a lot of interest and debate in medicine^{1–3}. They can be viewed as a simulation of an active therapy within a psychosocial context^{1–3}. Research in neurobiology has shown that placebo responses are accompanied by actual alterations in neural activity within brain regions involved in emotional regulation^{1–6}. Hence, rather than inducing a simple bias in response, placebos can induce actual biological effects and improve clinical outcomes. Among the cognitive and emotional factors that have been proposed to contribute to placebo effects, the interaction between the desire for symptom change and the expected symptom intensity has been proposed as a key component giving rise to placebo effects¹. In the psychology literature, this interaction is known as the desire-expectation model of emotions^{1,7–9}, which postulates that ratings of positive and negative emotional feelings are predicted by multiplicative interactions between ratings of desire and expectation. A number of experimental studies of placebo analgesia^{1,10,11} have corroborated the role of the desire-expectation model as a trigger of placebo effects. These findings have important implications for both clinical practice and clinical trials. On one hand, clinicians should harness the placebo effect to improve the clinical outcome of their patients (by managing expectations and desires through ethical use of suggestions and optimum caregiver-patient interactions)¹. On the other hand, assessment of expectation and desire levels is also important in clinical trials since placebo effects might strongly influence the results of a study. In unblinded trials, it is widely recognized that the overall effect attributed to a treatment might actually correspond to a combination of treatment and placebo effects. However, placebo effects might still play a role in blinded trials as well¹. For instance, blinded studies evaluating the effectiveness of acupuncture¹² and of implantation of human embryonic dopamine neurons into the brains of persons with severe Parkinson disease¹³ have shown that perceived treatment (or the treatment the participants thought they had received) can have stronger effects than the treatment actually received by the participants. These findings illustrate the relevance of measuring expectation, desire, and emotional levels in order to assess the contribution of placebo effects, and suggest that it is important to adjust for these variables when estimating treatment effects and interpreting the results of clinical trials¹. However, because it is generally impossible to rule out the presence of unmeasured confounders, simply measuring and adjusting for variables associated with placebo effects might not be enough to ensure a reliable estimation of the treatment effect. For instance,

Sage Bionetworks, Seattle, Washington, USA. Correspondence and requests for materials should be addressed to E.C.N. (email: elias.chaibub.neto@sagebase.org)

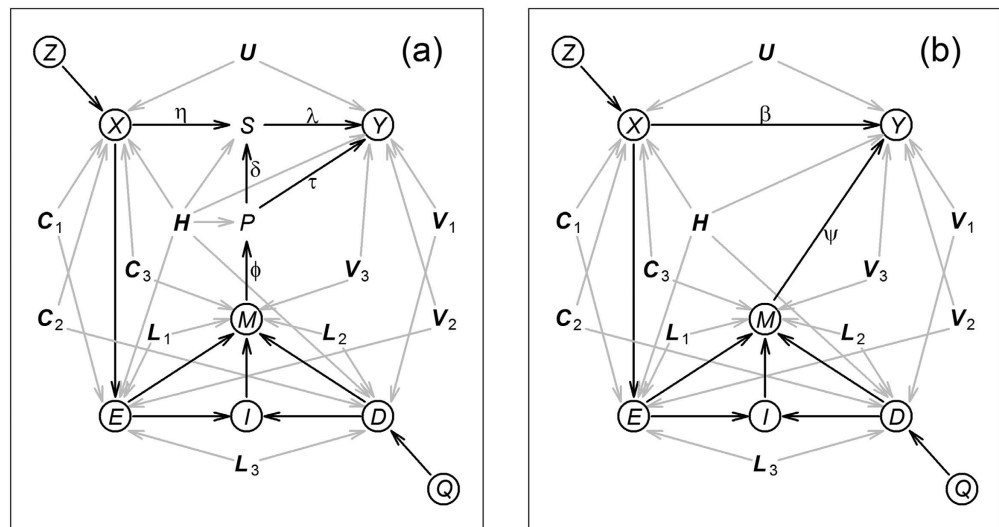


Figure 1. Direct acyclic graph representation of the causal model underlying the proposed IV approach for disentangling treatment and placebo effects in unblinded clinical trials. Circled and un-circled nodes represent observed and unobserved variables, respectively. Arrows represent the causal influence of a variable on another, with the influence of confounders on variables shown as grey arrows. The Z and X nodes represent, respectively, the participant's assigned and received treatment, whereas Q stands for the psychological encouragement treatment. The S and P variables represent the (unobserved) somatic and psychosomatic states of the participant, respectively. The E , D , I , and M nodes stand for the participant's expectation of symptom intensity, desire for improved symptoms, desire-expectation interaction, and emotional level, respectively. The sets of variables U , C_1 , C_2 , C_3 , L_1 , L_2 , L_3 , V_1 , V_2 , V_3 , and \mathcal{H} stand for unmeasured confounder variables. The Y node represents the outcome variable. Panel a shows the full model where the unobserved somatic and psychosomatic states of a participant are not directly represented in the causal model.

estimation based on regression models adjusting for the placebo related measurements still leads to biased estimates of the treatment effect, unless all confounders influencing the outcome variable enter the regression model.

The statistical method

Here we present a statistical approach to disentangle treatment and placebo effects using instrumental variables^{14–16} in randomized experiments. An instrumental variable (IV) is statistically independent from any unmeasured confounders, but is associated with the treatment variable and with the outcome variable (via its influence on the treatment variable alone). Use of IVs in randomized experiments allows the consistent estimation of treatment effects without the need to explicitly model the confounders (the technique even accounts for confounders the researcher is unaware about).

Our proposed method requires the ability to assess variables associated with placebo effects (e.g., levels of expectancy, desire, and emotion), and uses randomization to separately manipulate a pair of variables. The first, corresponds to a psychological encouragement variable aiming to increase the desire for improved symptoms. The study participants are randomized according to whether they receive the psychological encouragement or not. This “psychological treatment” IV allows the consistent estimation of the placebo effect on the outcome in the presence of confounders. The second, corresponds to a treatment assignment variable representing the random assignment of participants to active treatment or control therapy groups. It allows the estimation of the treatment effect on the outcome, after adjustment for the placebo effect. Mechanistically, the approach corresponds to a two-step procedure, which first estimates the contribution of the placebo effect on the outcome, and then the effect of the treatment on the residuals of the outcome variable after the contribution of the placebo effect has been removed.

A graphical representation of the causal model underlying our approach is given in Fig. 1a. Circled and un-circled nodes represent observed and unobserved variables, respectively. Arrows represent the causal influence of a variable on another, with the influence of unmeasured confounders shown as grey arrows. The binary variable Z represents the randomized treatment assigned to the participant (1 if participant is assigned to the active treatment group, and 0 if assigned to the control group), while X represents the treatment actually received by the study participant (1 if the participant receives the active treatment, and 0 otherwise). It is important to model both assigned and received treatment variables since participants won't necessarily subscribe to their assigned treatment, and the experiment might suffer from imperfect compliance.

The variable S represents the unmeasured biochemical/physiological (somatic) state of a participant and mediates the effect of the treatment on the outcome variable, Y . For instance, if X represents a drug treatment, then S could represent the physiological state induced by the biochemical pathways targeted by the drug. The causal effects of X on S and of S on Y are quantified, respectively, by η and λ . The outcome variable is also influenced by the unmeasured psychosomatic state of the participant, represented by P . We allow P to influence Y via a direct

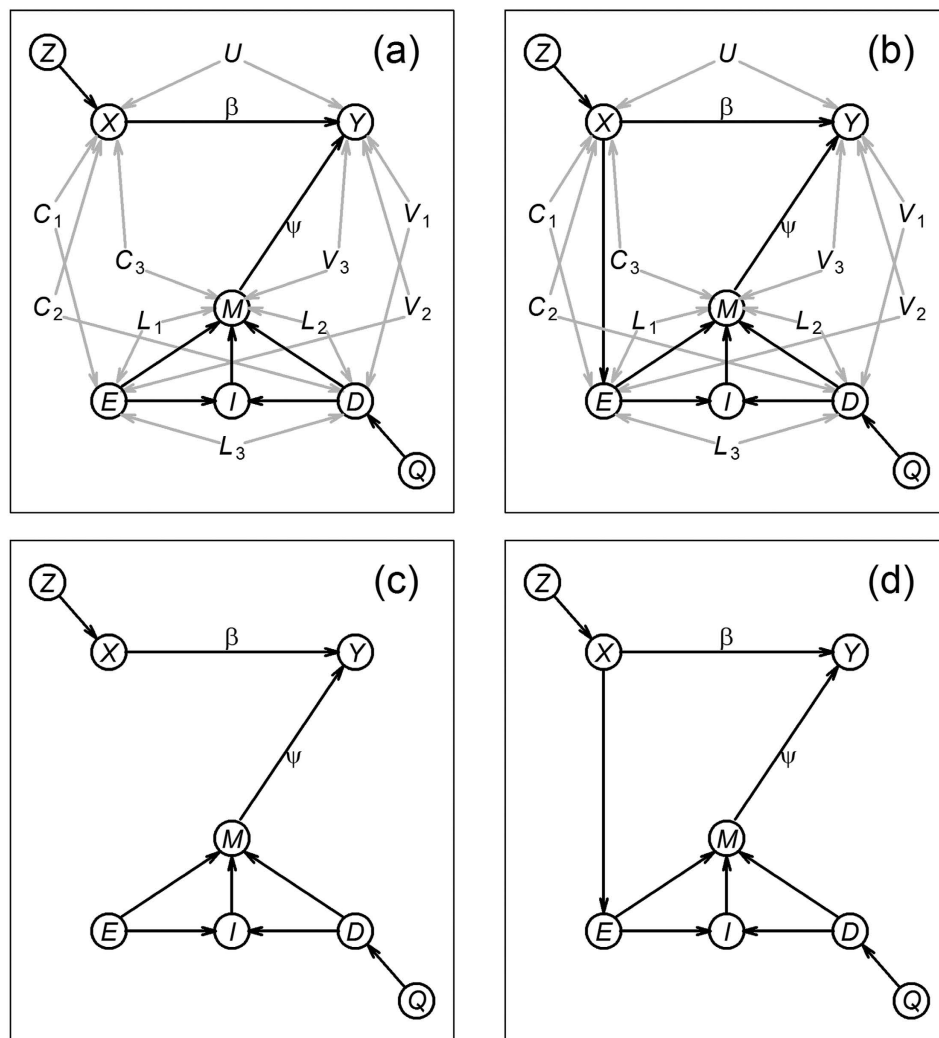


Figure 2. Models used in the simulation study. Node definitions are provided in Fig. 1. Panels a and b represent, respectively, blinded and unblinded trials influenced by confounders. For simplicity we include a single confounder variable per pair of endogenous variables (other than I), but still simulate confounding across the 10 possible pairwise combinations of the endogenous variables X , Y , E , M , and D . Panels c and d represent, respectively, unconfounded blinded and unblinded trials. For simulations under the null $H_0: \psi = 0$ there are no arrows from M to Y . Similarly, for simulations under $H_0: \beta = 0$, there are no arrows from X to Y .

path, quantified by τ , and by an indirect path, mediated by S , and quantified by the product $\delta\lambda$. The combined effect of the direct and indirect paths represents the placebo effect. The direct path from P to Y represents the influence of the psychosomatic state on the outcome mediated by biochemical and physiological pathways distinct from the pathways influenced by the active treatment, while the influence of P on S allows for the possibility that P also influences the same pathways targeted by the treatment X . (Experimental evidence that placebo effects influence biochemical pathways is provided, for example, in studies of placebo analgesia involving endogenous opioid systems^{1,17–22}. See also Fig. 2 in ref. 2, for empirical support about pathways influenced by both psychosocial context and drug treatments).

The role played by the expectation-desire model of emotions is made explicit by the observed variables E , D , I and M , representing, respectively, the expected symptom intensity, the desire for symptom improvement, the interaction between expectation and desire, and the emotional level (measured, for example, by the participant's mood). According to the expectation-desire model, M is directly influenced by E , D , and their interaction $I = E \times D$. The causal influence of M on P is quantified by ϕ .

In unblinded trials it is reasonable to expect that the treatment actually received by the participant will affect its expected symptom intensity, since participants who know they are receiving the active treatment will more likely experience an increase in their expectation to feel better. Hence, we include an arrow from X to E . The implication is that the treatment can influence the outcome not only via the participant's somatic state, but also by its psychosomatic state via the paths $X \rightarrow E \rightarrow M \rightarrow P$ and $X \rightarrow E \rightarrow I \rightarrow M \rightarrow P$. The binary variable Q represents the

randomized psychological encouragement IV assuming the value 1 when a encouragement conversation/interaction (aiming to increase the desire for symptom improvement) is applied to the participant, and 0 otherwise.

In addition to the key variables described so far, it is important to recognize the existence of unmeasured confounders. Except for the exogenous variables Z and Q , that by construction are not associated with any unmeasured confounders, the model includes confounders influencing all pairs of endogenous variables other than I , namely, X, E, D, M, P, S , and Y . (It is not necessary to include confounders between I and the other endogenous variables, since I is deterministically computed as the product of E and D). For instance, U represents a set of unmeasured confounder variables influencing X and Y . In order to avoid cluttering the figure, the confounder variables influencing S and P and all other endogenous variables are represented by the vector of variables $\mathcal{H} = (\mathbf{H}_1, \dots, \mathbf{H}_1)^T$. (For the same reason the figure does not explicitly shows the error terms, which account for unmeasured variables influencing each particular variable in the model and are uncorrelated with each other). It would be unrealistic to assume, for example, that the emotion of a participant is determined by E, D , and I alone. Hence, the model allows sets of unmeasured confounders, such as L_1, L_2 and L_3 , to influence emotion and expectation, emotion and desire, and expectation and desire, respectively. Similarly, it would be unrealistic to assume that emotion alone influences the psychosomatic state of a participant, and the model accommodates unmeasured confounders influencing these variables as well. Although, in practice, not all endogenous variables (other than I) will necessarily be influenced by confounders, the model still includes confounders for all 21 pairwise combinations of endogenous variables, since we want to derive estimators for the placebo and treatment effects under the most general setting possible.

In practice, however, it is impossible to accurately measure the unobserved somatic and psychosomatic states of a participant. Hence, Fig. 1b presents a reduced version where S and P are not explicitly represented in the graph. Assuming linear relationships between S and X, P and M , and Y, S , and P , the causal influence of X on Y is given by $\beta = \eta\lambda$, while the influence of M on Y is given by $\psi = \phi\tau + \phi\delta\lambda$. Under this reduced model the instrumental variable Q allows for the consistent estimation of the net placebo effect, ψ , using the IV estimator $\hat{\psi}_{IV} = \widehat{\text{Cov}}(Q, Y)/\widehat{\text{Cov}}(Q, M)$. Once the net placebo effect is estimated, it is possible to estimate the causal effect of X on Y using the IV estimator of the causal effect of X on the residuals of the outcome variable after the removal of the placebo effect, $\hat{\beta}_{2sIV} = \widehat{\text{Cov}}(Z, \hat{R})/\widehat{\text{Cov}}(Z, X)$, where $\hat{R} = Y - \hat{\psi}_{IV}M$ (see Methods for details).

Performance evaluation

We assessed the statistical performance of the proposed method (and compare it to a naive regression approach) in 16 simulation experiments evaluating the empirical type I error rate and empirical power of randomization tests for the null hypotheses that the placebo effect is zero, $H_0: \psi = 0$, and that the treatment effect is zero, $H_0: \beta = 0$. Descriptions of the randomization tests and simulation experiments are provided in the Methods. We simulated data from blinded and unblinded trials, in the presence and absence of confounders, according to the models presented in Fig. 2.

For each setting, we ran 4 separate simulation experiments generating data: (i) under the null for treatment and placebo effects; (ii) under the alternative for treatment, and null for placebo effects; (iii) the other way around; and (iv) under the alternative for treatment and placebo effects. Each simulation experiment employed 10,000 distinct synthetic data sets with diverse characteristics (see Methods). Although the randomization tests are non-parametric procedures free of distributional assumptions, we still generated data using gaussian errors in order to meet the distributional requirements of the regression based analytical tests used in our comparisons.

Figure 3 presents the results for the placebo effect tests, and shows that the error rates of the IV approach (red and blue) are controlled at the exact nominal levels in both blinded and unblinded settings, in the presence and absence of confounders. The regression approach (brown and dark-orange), on the other hand, shows highly inflated errors in the presence of confounders (Fig. 3a and b), since the association between M and Y , caused by confounders, is mistaken by an influence of M on Y . Being able to control type I error rates at the exact nominal level is a desirable statistical property, as it means that the test is neither conservative nor liberal.

Figure 4 presents the results for the treatment effect tests in the blinded setting. In addition to the two-step estimator (blue), we also evaluated the simple IV estimator $\hat{\beta}_{IV} = \widehat{\text{Cov}}(Z, Y)/\widehat{\text{Cov}}(Z, X)$, which does not account for the placebo effect (red). The results show, again, well controlled error rates for both IV approaches, but inflated errors for the regression test (brown) in the presence of confounders (Fig. 4a and b).

Figure 5 presents the results for the unblinded case. All panels show slightly inflated errors for the two-step IV estimator (blue). The likely reason is that the estimated placebo effects are noisy and unable to completely block the influence of X on Y through the paths mediated by M . To test this supposition, we evaluated an additional IV estimator, where the true placebo effect was used in the computation of the residuals (i.e., we estimated β by $\hat{\beta} = \widehat{\text{Cov}}(Z, R)/\widehat{\text{Cov}}(Z, X)$, where $R = Y - \psi M$, instead of $\hat{\beta}_{2sIV} = \widehat{\text{Cov}}(Z, \hat{R})/\widehat{\text{Cov}}(Z, X)$, where $\hat{R} = Y - \hat{\psi}_{IV}M$). Results based on this estimator (dark-orange) show that, indeed, adjustment by the true placebo effect leads to error rates controlled at the nominal level. The regression approach (brown) shows again highly inflated errors in the presence of confounders (Fig. 5a and b). Figure 5a and c show well controlled errors for the non-adjusted IV estimator (red) in the absence of placebo effects as, in this case, there are no paths from X to Y , and the association between X and Y induced by confounders is accounted by the IV estimator. Figure 5b and d, on the other hand, show highly inflated error rates in the presence of placebo effects since, in this case, X can influence Y through the paths mediated by M .

These observations suggest that, in practice, when analyzing the results of unblinded trials, we should first test for the existence of placebo effect, and then use the two-step IV estimator if $H_0: \psi = 0$ is rejected, and the non-adjusted one if $H_0: \psi = 0$ is accepted. While this strategy can decrease the chance of the two-step approach making a type I error in the absence of placebo effects, the estimator is still unable to avoid slightly inflated errors produced in the presence of placebo effects. We point out, however, that the two-step procedure still represents

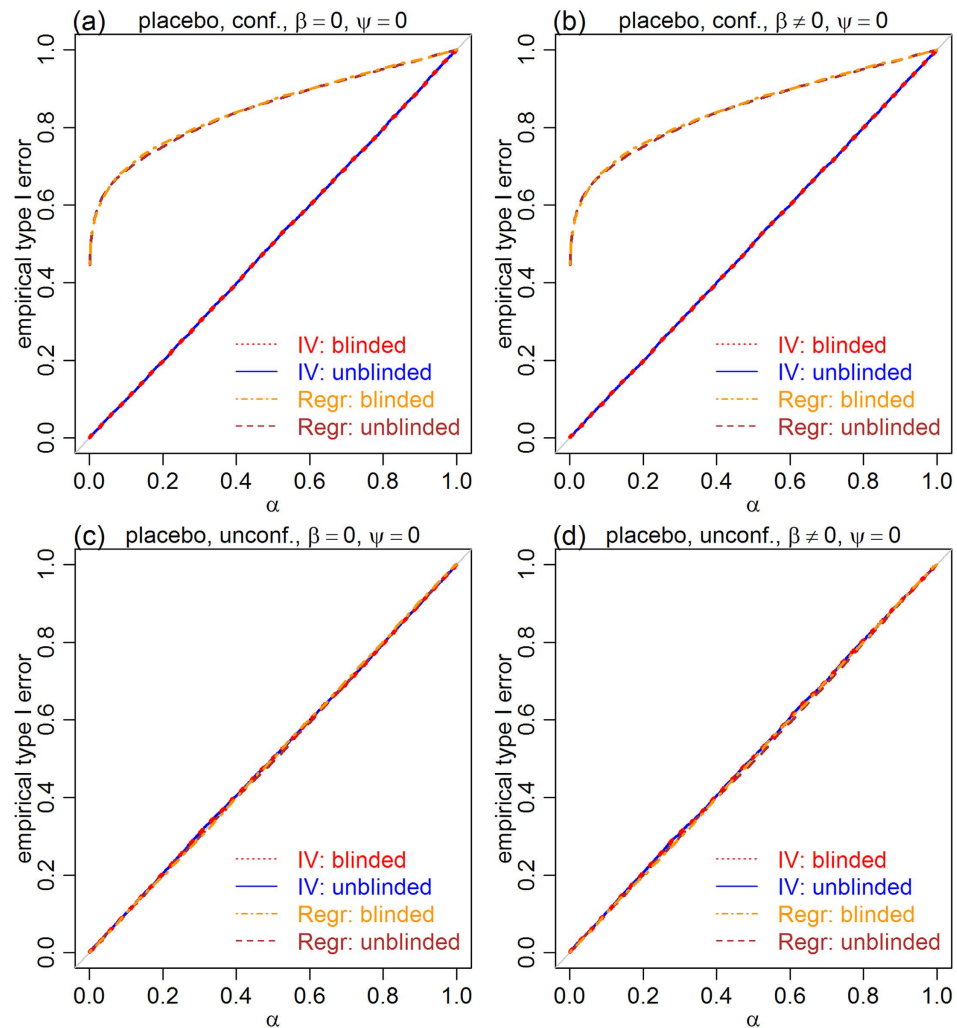


Figure 3. Empirical type I error rates of the placebo effect null, $H_0: \psi = 0$, in both blinded and unblinded settings. Panels a and b show that, in the presence of confounders, the type I error rate of the IV approach is controlled at the exact nominal level (red and blue), whereas the regression based test leads to highly inflated error rates (orange and brown). Panels c and d show that, in the absence of confounding, both IV and regression approaches show well controlled errors. The nominal significance level is represented by α .

a strong improvement over the alternative approach of not adjusting for placebo effects in the presence of confounders (compare the red and blue curves in Fig. 5b).

For completeness, we also report an evaluation of the empirical power (Suppl. Figs 1, 2 and 3). We point out, however, that power results are more sensitive to the choice of parameter values employed in the generation of the simulated data (e.g., sample size, the strength of treatment, placebo and confounding effects, and etc), than the type I error rates. In any case, these empirical power results, still serve to illustrate some general patterns. For instance, the regression tests tended to show considerably stronger power than the IV approaches in the presence of confounders (compare the brown and blue curves in panels a and b of Suppl. Figs 2 and 3). We point out, however, that this increased power is likely an artifact of the biased estimates of β outputted by the regression approach. Suppl. Fig. 4, illustrates how the regression estimates tended to show larger bias than the estimates generated by the IV estimators (note the heavier tails of the brown density, in both blinded and unblinded cases). In other words, the increased power is likely a consequence of the overestimation of the treatment effect by the regression approach, which mistakenly interprets the association between treatment and outcome caused by unmeasured confounders as a stronger influence of the treatment on the outcome.

At least for the parameter ranges adopted in our simulations, we observed good empirical power of the IV approach to detect placebo effects, even when the correlation between psychological encouragement and emotional level was relatively low (Suppl. Fig. 5a). This suggests that the psychological encouragement treatment does not need to be highly effective in manipulating the emotional levels, in order for the approach to work well in practice. Similarly, Suppl. Fig. 5b shows good empirical power of the two-step IV approach to detect treatment effects when the correlation between the assigned and received treatment is moderate, suggesting that the proposed approach does not require high levels of compliance in order to perform well.

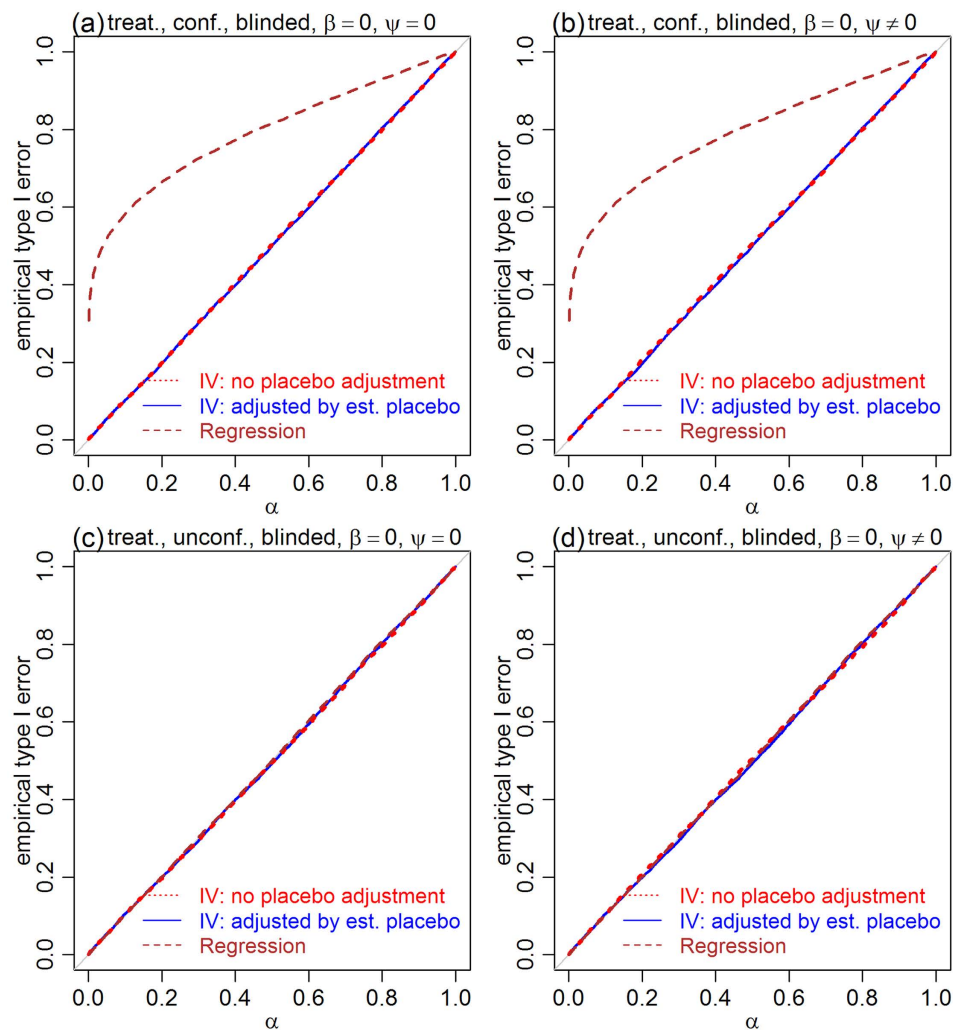


Figure 4. Empirical type I error rates for the treatment effect null, $H_0: \beta = 0$, in the blinded setting. Panels a and b show that, in the presence of confounders, the type I error rates of the IV approaches are controlled at the exact nominal level (red and blue), whereas the regression based test leads to highly inflated error rates (brown). Panels c and d show that, in the absence of confounding, both IV and regression approaches show well controlled errors. The nominal significance level is represented by α .

A natural question, at this point, is whether larger sample sizes (and, hence, more precise estimates of ψ) would be able to decrease the slightly inflated error rates produced by the two-step estimator in unblinded trials. Figure 6 presents additional simulation experiments showing that, while the empirical power and the $\hat{\psi}_{IV}$ and $\hat{\beta}_{2sIV}$ estimates are greatly improved by larger sample sizes, the type I error rates stay roughly the same (likely because larger sample sizes increase the ability of a test to detect small effects, since the randomization null distributions tend to be more concentrated around 0, so that the improved $\hat{\psi}_{IV}$ estimates are counterbalanced by the increased propensity to detect small and spurious treatment effects). These results suggest that special care must be taken while interpreting the results of hypothesis tests in the unblinded case, even for large sample sizes. In any case, when the goal is estimation rather than testing, the consistency of the two-step estimator guarantees that the treatment estimates will converge to the true value as the sample size increases.

This observation is particularly important in view of the current trend in the biomedical field, where researchers are shifting from relying exclusively in p-values and are paying more attention to parameter estimates and confidence intervals. To meet this latter need, we also describe in the Methods how to generate confidence intervals (CIs) for placebo and treatment effects by inverting randomization tests. Figure 7 shows 95% CIs for the placebo and treatment effects, from 3 simulated data sets of increasing sizes. The randomization CIs inherit the statistical properties of the randomization tests, hence, the placebo effect CIs (and treatment effect CIs from blinded trials) are exact in the sense that a $100(1 - \alpha)\%$ interval will contain the true parameter value $100(1 - \alpha)\%$ of the time. Note that while the treatment effect CIs from unblinded trials won't be exact, they are still going to be centered around the estimated treatment effect, which will, nevertheless, converge to the true value as the sample size increases.

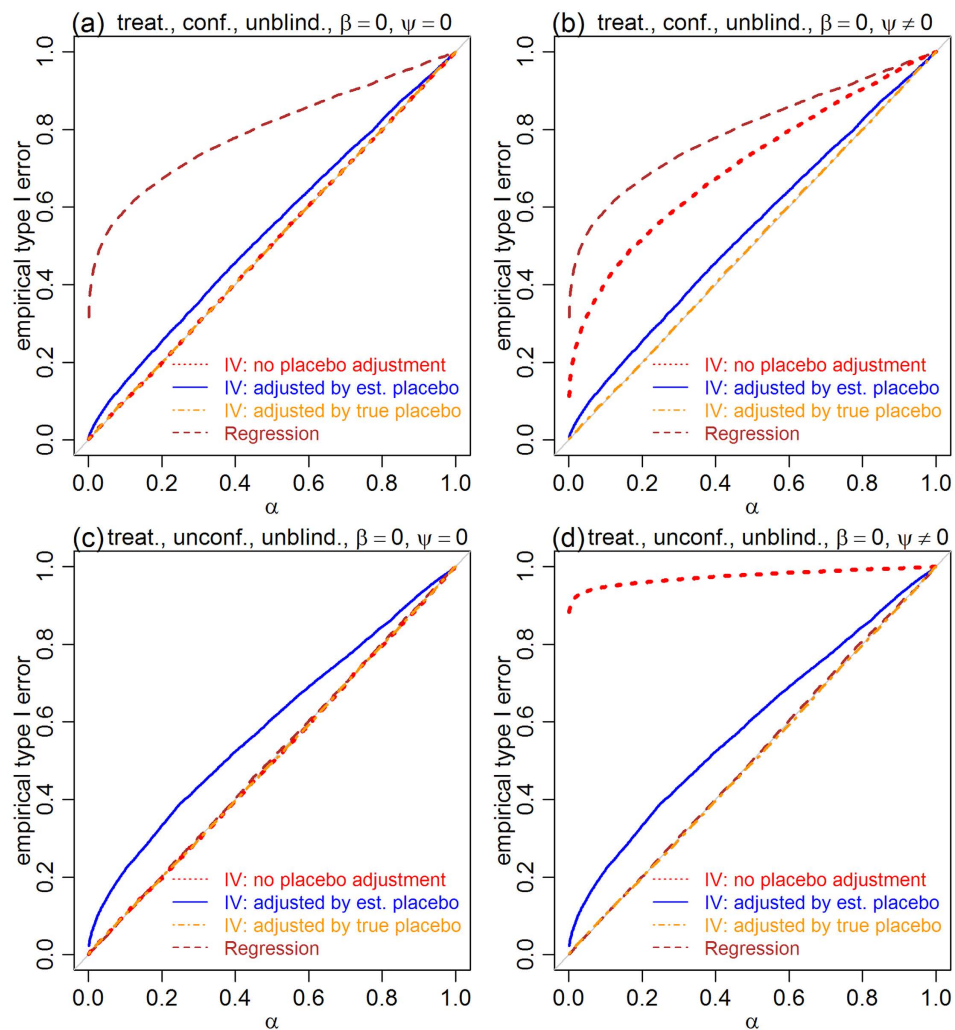


Figure 5. Empirical type I error rates for the treatment effect null, $H_0: \beta = 0$, in the unblinded setting. The two-step IV approach (blue) shows slightly inflated errors in the presence (panels a and b) and absence (panels c and d) of confounders. Note that the larger errors in panels c and d, in comparison to a and b, are likely due to the effective stronger influence of X on M in the simulations unaffected by confounders (the presence of confounders can considerably increase the amount of noise), so that adjustment by $\hat{\psi}_{IV}$ leaks more information about X in the absence than in the presence of confounders. The estimator adjusted by the true placebo effect (dark-orange) leads, nonetheless, to well controlled errors. The non-adjusted IV approach (red) leads to well controlled errors in the absence of placebo effects (panels a and c), but to highly inflated errors in the presence of placebo effects (panels b and d). Regression (brown) leads to highly inflated errors in the presence of confounders (panels a and b), but to well controlled error rates in their absence (panels c and d).

Discussion

Clinical trials traditionally employ blinding to control the influence of placebo effects. It has been pointed out, however, that even blinded studies might be influenced by placebo effects, as the patients' perceptions and beliefs about the treatment they think they received are able to trigger strong placebo effects^{1,12,13}. Recently, a number of statistical approaches have been proposed to quantify the contributions of treatment and placebo effects to a clinical outcome^{23–25}. These approaches, nonetheless, are tailored to blinded trials, and leverage blinding assessment data to quantify the amount of unmasking taking place during the trial. Our IV approach, on the other hand, actively manipulates emotion levels and allows the quantification of treatment and placebo effects not only in blinded, but also in unblinded trials.

The key idea underlying the IV approach (what actually allows the consistent estimation of both treatment and placebo effects in the presence of unmeasured confounders), is the use of randomization to separately manipulate the treatment assignment and encouragement conversations/interactions. In this sense, the proposed approach is similar in spirit (but not exactly equivalent) to a randomized treatment-belief trial (RTB)²⁶, where the treatment assignment is manipulated by randomization, whereas the belief is manipulated by varying the allocation ratio of participants assigned to control and treatment groups in a, necessarily, blinded trial. Hence, our IV approach can be viewed as a more flexible type of RTB that is applicable to both blinded and unblinded studies, and might be

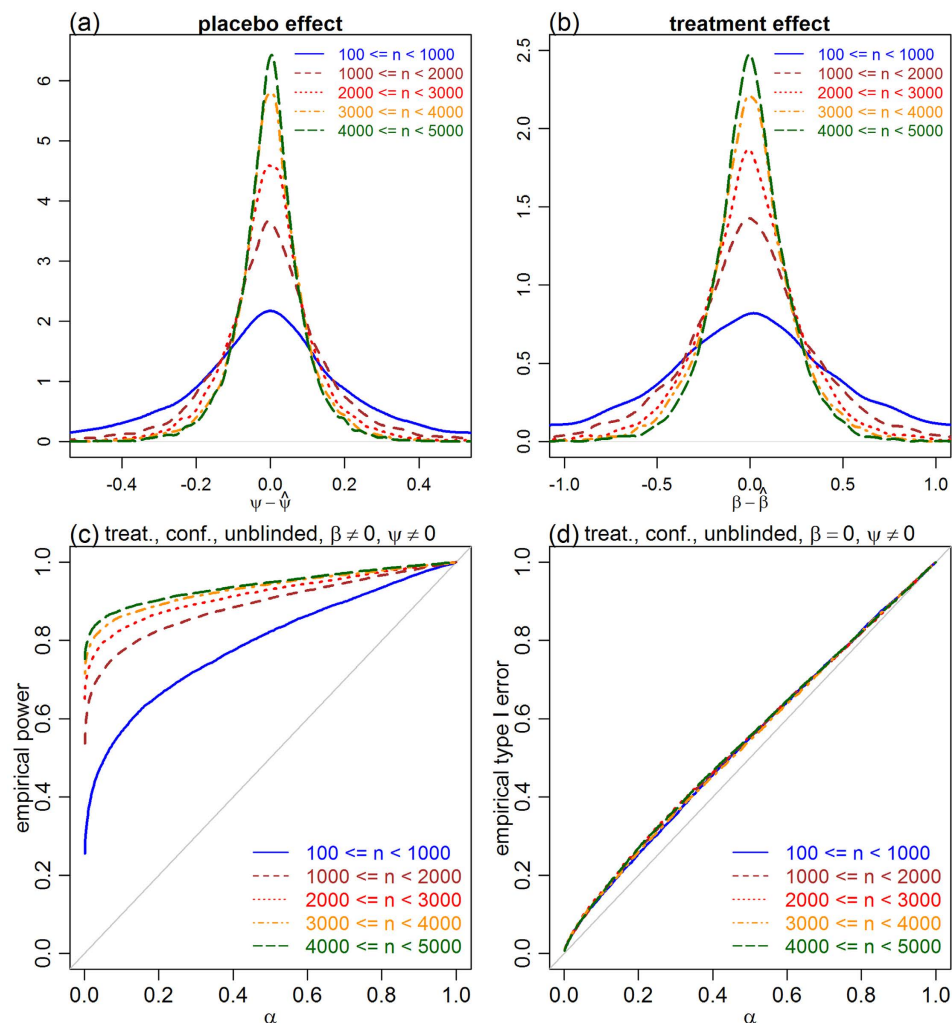


Figure 6. Consistency of the $\hat{\psi}_{IV}$ and $\hat{\beta}_{2sIV}$ estimators. Panels a and b present, respectively, the densities of $\psi - \hat{\psi}_{IV}$ and $\beta - \hat{\beta}_{2sIV}$ for 5 increasing sample size ranges, and illustrate the consistency of the $\hat{\psi}_{IV}$ and $\hat{\beta}_{2sIV}$ estimators (which tend to get closer to the true parameter values as the sample size increases). Panel c shows that, as expected, the statistical power to detect a treatment effect increases with the sample size. Panel d, on the other hand, shows that increasing sample sizes do not reduce type I error rates, even though we are able to better estimate the placebo effects. The likely reason is that while larger sample sizes lead to better $\hat{\psi}_{IV}$ estimates, they also increase the statistical power to detect very small effects, so that the advantage of a more precise estimate of $\hat{\psi}_{IV}$ is counterbalanced by the increased propensity to detect small and spurious treatment effects as true signals. Results were based on data simulated from unblinded trials influenced by placebo effects and confounders, as described in the Methods section.

easier to administer than a standard RTB, which requires the stratification of study participants over several arms with distinct treatment/control allocation ratios in order to be able to assess placebo effects.

The proposed IV approach enjoys appealing statistical properties. The IV estimators are consistent, meaning that the estimates converge to the true values as sample size increases. The randomization tests for placebo effects are exact in both blinded and unblinded trials, whereas the treatment effect tests are exact in blinded trials, but slightly liberal in unblinded ones. Furthermore, the confidence intervals obtained by inverting randomization tests inherit these appealing properties.

Under our proposed statistical model, the identification of the average treatment effect, β , and the average placebo effect, ψ , requires a number of assumptions, namely: (i) that $Z(Q)$ do not share common causes with Y ; (ii) that $Z(Q)$ is marginally associated with $X(M)$; (iii) that Z influences Y only through X (and Q influences Y only through M); (iv) that Y is linearly associated with X and M ; (v) that the influence of residual errors and unmeasured confounders on Y is additive; and (vi) that the average causal effects, β and ψ , are constant across all individuals in the population.

The first 3 assumptions correspond to the “core” conditions for instrumental variables¹⁶. They allow the identification of bounds²⁷ for the causal effect (i.e, lower and upper limits for the effect that are consistent with the data), but are not sufficient to identify a point estimate for the causal effect. In observational studies, conditions i to iii need to be carefully evaluated in order to assess the validity of the putative instrument. However, in the

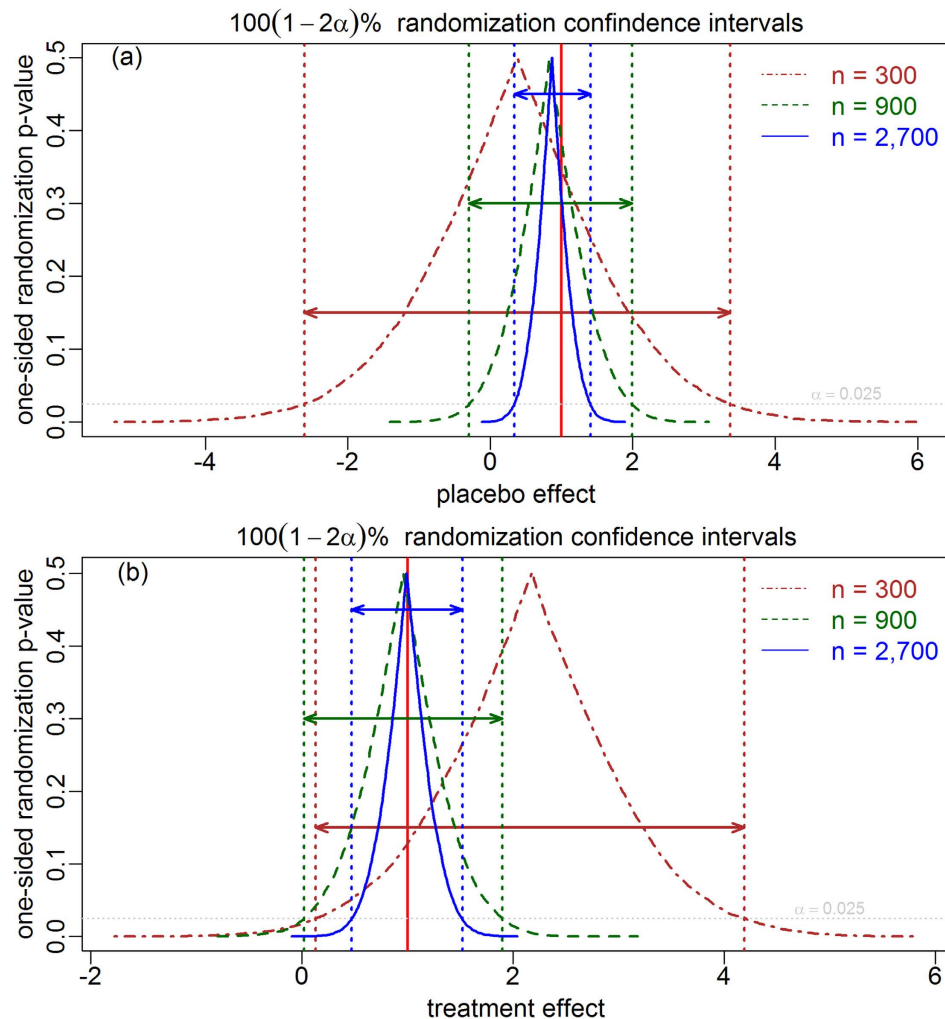


Figure 7. Randomization confidence intervals for placebo and treatment effects. The brown, dark-green and blue curves show the one-sided p-value profiles derived from randomization tests for 3 simulated data sets of increasing sizes (300, 900, and 2,700, respectively), generated under the unblinded setting influenced by confounders (all simulation parameters, other than sample size, were set to 1). The 95% confidence intervals for the placebo (panel a) and treatment effects (panel b) are shown by the respective double-headed colored arrows. The red vertical line corresponds to the true parameter values, $\psi = 1$ and $\beta = 1$.

context of our proposed approach, where the instruments are randomized, condition *i* is valid by construction, while condition *ii* holds if there is some degree of compliance between the randomly assigned treatment and the treatment effectively adopted by the study participants, and if the psychological suggestions are able to manipulate the emotion level of the study participants (observe, however, that these assumptions can be checked empirically, by inspecting $\widehat{\text{cor}}(Z, X)$ and $\widehat{\text{cor}}(Q, M)$). Condition *iii*, also known as the exclusion restriction, is only guaranteed to hold in double-blinded trials²⁸, since knowing to which treatment arm a participant has been allocated might change the participant's behavior in ways that affect the outcome other than through the treatment and/or placebo effects. For instance, condition *iii* would be violated if assignment to the treatment group increased awareness about risk factors, and the participants adopted preventive measures that they would not have adopted, had they been assigned to the control group²⁹. In any case, condition *iii* is often reasonable in other experimental situations (but still needs to be stated as an assumption).

The additional assumptions *iv* to *vi* allow the identification of point estimates for the causal effects and are specific to the structural equations underlying our proposed method. The adequacy of assumptions *iv* and *v* can, nonetheless, be checked empirically by examining if a linear model provides a reasonable fit to the data. Assumption *vi*, on the other hand, is often times more contentious since it is unlikely to (strictly) hold in most applications based on continuous responses, and is generally impossible to hold for binary responses³⁰. (An alternative causal framework based on potential outcomes^{31,32} explicitly allows for effect heterogeneity by focusing on unit-level participant specific causal effects, but at the expense of only being able to identify the treatment effect for a latent subpopulation of “compliers”³³).

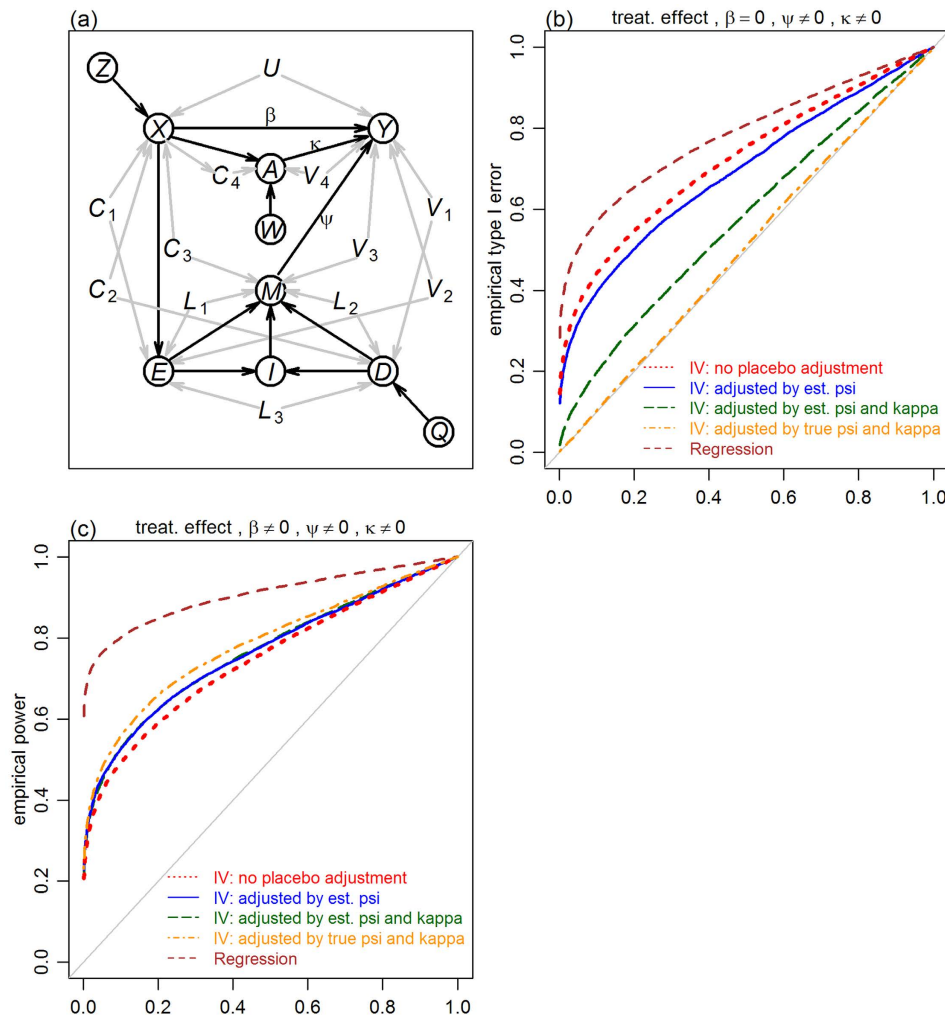


Figure 8. A more complex example. Panel a presents a more complex model where the placebo effect is mediated by M (according to the desire-expectation model) but also by an additional variable A (definitions of the additional nodes in the graph are provided in Fig. 1). Assuming that a randomized instrument, W , is available to manipulate A , we can estimate the treatment effect using the estimator $\hat{\beta}^* = \widehat{\text{Cov}}(Z, \hat{R}^*) / \widehat{\text{Cov}}(Z, X)$ where $\hat{R}^* = Y - \hat{\psi}_{IV}M - \hat{\kappa}_{IV}A$. Panel b shows the empirical type I error rates for a simulation experiment under the unblinded setting influenced by confounders. The IV estimator adjusted by the true ψ and κ values is able to control error rates at the nominal levels (dark-orange). The IV estimator adjusted by $\hat{\psi}_{IV}$ and $\hat{\kappa}_{IV}$ shows slightly inflated errors (dark-green). As expected, adjustment with $\hat{\psi}_{IV}$ alone (blue) leads to higher error rates than adjustment with both $\hat{\psi}_{IV}$ and $\hat{\kappa}_{IV}$. Similarly, the IV estimator using no adjustment (red) has higher errors than adjustment by $\hat{\psi}_{IV}$ alone. The regression based estimator (brown) is adjusted by both M and A covariates, but still leads to inflated errors due to the presence of confounders. Panel c shows the empirical power results.

An implicit assumption of the model in Fig. 1a is that the placebo effect is mediated exclusively by the interplay of (perfectly measured levels of) desire, expectation, and emotion, assessed at a single time point. While it is believed that the desire-expectation model plays a key role in the triggering of placebo effects, other mechanisms, such as conditioning and learning, might also be at work¹⁻³. Clearly, when this is the case, a treatment effect estimate, adjusted by the desire-expectation component alone, will still be biased (although less biased than the estimate computed without accounting for it). In any case, if we are also able to assess and measure these additional mechanisms, then the same statistical framework can be used to obtain consistent estimates of treatment effects in the presence of confounders (we only need additional IVs to manipulate the additional placebo related variables). Figure 8 shows an example.

All simulation results presented in the main text were generated under the assumption of perfect measurements of the emotion level variable. In practice, however, recording of emotion levels might be more prone to the influence of measurement error than the recording of the treatment, outcome and instrumental variables. In order to evaluate the influence of measurement error in the performance of the proposed method, we conducted a number of additional simulation experiments with data generated under the presence of varying amounts of measurement error over the emotion level variable. A detailed description of the simulation design and results is presented in the Supplementary Note. Our results suggest that, at least for the settings evaluated in our

simulations, the placebo effect IV estimator, $\hat{\psi}_{IV}$, tended to be resilient to the effects of measurement error, and was considerably less biased than the regression approach estimator, $\hat{\psi}_R$. This result is not surprising given that the initial motivation for the use of IV estimators in economics was to handle measurement error in explanatory variables¹⁴, and, hence, one would expect $\hat{\psi}_{IV}$ to be able to account for measurement error on the emotion level variable. Furthermore, our results also illustrated that the placebo adjusted IV estimator of treatment effects, $\hat{\beta}_{2sIV} = \widehat{\text{Cov}}(Z, \hat{R}) / \widehat{\text{Cov}}(Z, X)$, tended to be less biased than the regression, $\hat{\beta}_R$, and unadjusted IV estimators, $\hat{\beta}_{IV} = \widehat{\text{Cov}}(Z, Y) / \widehat{\text{Cov}}(Z, X)$, in the presence of measurement error, although the decrease in bias achieved by $\hat{\beta}_{2sIV}$ tended to be less accentuated in comparison to the decrease observed for the placebo effect. This last observation is also not surprising since the estimation of the placebo effect is never free from noise, and, even though $\hat{\psi}_{IV}$ seems to be able to reduce the additional bias induced by measurement error, it cannot completely neutralize it. Hence, in the presence of measurement error, the placebo effect estimates, $\hat{\psi}_{IV}$, employed in the computation of the residuals, \hat{R} (which enter the estimation of $\hat{\beta}_{2sIV}$), tend to be less effective in removing the influence of the placebo effect on the outcome variable.

The current popularity of IV methods in observational studies seems to arise from their ability to account for unmeasured confounding. However, an increasing body of literature shows that IV methods can be very sensitive to violations of the underlying assumptions. Well known sources of biases in IV analysis include: bias amplification due to weak association between instrument and exposure/endogenous explanatory variable²⁸; violations of the exclusion restriction²⁸; and biases generated by selection mechanisms^{34,35}. It has also been argued in the literature^{36,37} that IV methods shift the problem of measuring and adjusting for confounders of the treatment-outcome association, to the problem of dealing with confounders of the instrument-outcome association. We point out, nonetheless, that because our proposed approach is based on truly randomized IVs, it avoids instrument-outcome confounding issues. However, our method is still vulnerable to bias amplification, to selection bias issues (including selection of treatment³⁵ in situations where an analyst focus on only two treatments while ignoring other possibilities, e.g., no treatment), and to violations of the exclusion restriction.

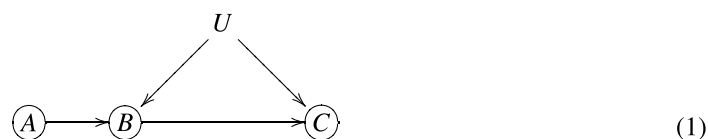
From a pragmatic perspective, the proposed method is (in principle) easy to implement. It only requires the ability to assess expectation, desire, and emotion, as well as, the development of a psychological encouragement IV, capable of manipulating the level of desire of a study participant. For example, in trials run into a clinic, a simple encouragement conversation/interaction with a caregiver would work as the “active treatment” of the psychological encouragement IV. The desire and emotion level could then be recorded by a questionnaire or interview after the encouragement treatment, but prior to the measurement of the outcome variable.

Another application of the proposed method (the one that actually motivated this work) is in the personalized monitoring of treatment response in mobile health. The statistical validity of using treatment assignment as an IV, in the context of longitudinal data provided by a single patient, has been established in ref. 38. However, as pointed out by the authors, it is impossible to disentangle treatment and placebo effects based on the treatment assignment IV alone, since it is impossible to blind the patient to a self administered treatment. Implementation of the proposed IV approach in mobile health applications is also in principle strait-forward. For instance, the psychological treatment could be delivered by encouragement messages popping up in the screen of a smartphone (according to a randomized schedule, where every day the participant has an equal chance of receiving, or not, the encouragement message), and the measurement of the emotion and desire levels can be assessed by short electronic surveys/questionnaires delivered by the participant’s smartphone on a daily basis. We expect that the proposed method might play an important role in these personalized medicine^{39,40} applications.

Finally, for both (population-based) clinical trials and personalized monitoring of treatment response, the instrument Q serves the double role of disentangling placebo from treatment effects, and increasing the desire for improved symptoms. This latter capacity can induce a placebo effect and ultimately lead to more positive clinical outcomes. While the manipulation of the expectation for symptom intensity could, in principle, be used to consistently estimate a placebo effect under the proposed approach (i.e., we could have an IV influencing E instead of D), the manipulation of expectation levels needs to be accompanied by the honest disclosure of the expected benefits of a treatment (and, in some cases, might raise ethical issues)². Manipulation of the desire for improved symptoms, on the other hand, provides an ethically defensible practice in the design of clinical trials and in the personalized monitoring of patients.

Methods

Identification of causal effects using instrumental variables. We subscribe to the mechanism-based account of causation championed by Pearl⁴¹. In the mechanism-based framework, the qualitative description of the assumptions regarding the causal relations between the variables is encoded in a directed acyclic graph (DAG). When confounder variables are present, it is still possible to use instrumental variables to identify causal effects, whenever certain parametric and distributional assumptions hold. To fix ideas consider the toy causal model,



where A represents an IV, C represents the outcome variable, B represents an exposure/endogenous explanatory variable with a causal influence on the outcome C , and U represents an unmeasured confounder influencing both

B and C . Three necessary (although not sufficient) conditions¹⁶ required for A to qualify as an IV include that: (i) A must be marginally independent of any confounders of B and C , that is, $A \perp\!\!\!\perp U$; (ii) A must be marginally associated with B , that is, $A \not\perp\!\!\!\perp B$; and (iii) the association between A and C must be mediated exclusively by B , that is, $A \perp\!\!\!\perp C \mid (B, U)$. Inspection of Fig. 1a and b shows that these three assumptions are satisfied for the instrumental variable Q relative to the emotion level M and the outcome Y , as well as, for the instrumental variable Z relative to the received treatment X and outcome Y . In the context of randomized clinical trials, assumption *i* is valid by construction due to the randomization of the instruments. Assumption *ii* holds if there is some degree of compliance between the randomly assigned treatment and the treatment effectively adopted by the study participants (i.e., $Z \not\perp\!\!\!\perp X$), and if the psychological suggestions are able to manipulate the emotional level of the study participants (i.e., $Q \not\perp\!\!\!\perp M$). Assumption *iii* is only guaranteed to hold in double-blinded trials²⁸, but is often reasonable in other experimental situations.

As we describe next, the identification of the causal effects of M on Y and of X on Y requires, nonetheless, the additional assumptions of linear relations between Y and M and between Y and X . Assuming a linear relation between the outcome, Y , and the unobserved somatic and psychosomatic state variables, S and P , we have that,

$$Y = \mu_Y + \lambda S + \tau P + f_Y(\mathbf{U}, \mathbf{V}, \mathbf{H}) + \varepsilon_Y, \quad (2)$$

where $\mathbf{V} = (\mathbf{V}_1, \mathbf{V}_2, \mathbf{V}_3)^T$, $\mathbf{H} = (\mathbf{H}_1, \dots, \mathbf{H}_{11})^T$, ε_Y represents an error term accounting for the unmeasured variables influencing exclusively Y , and $f_Y(\mathbf{U}, \mathbf{V}, \mathbf{H})$ represents a general scalar function of the variables in $(\mathbf{U}, \mathbf{V}, \mathbf{H})$ influencing Y .

Since S and P are unobserved variables, we need to derive the reduced model for the outcome variable that is not a function of S and P . Assuming a linear relation between P and M , and between S and P and X , we have that,

$$P = \mu_P + \phi M + f_P(\mathbf{H}) + \varepsilon_P, \quad (3)$$

$$S = \mu_S + \eta X + \delta P + f_S(\mathbf{H}) + \varepsilon_S, \quad (4)$$

where $f_P(\mathbf{H})$ and $f_S(\mathbf{H})$ are arbitrary scalar functions of \mathbf{H} , and ε_P and ε_S are the respective error terms influencing P and S , respectively (we also assume that all variable specific error terms, ε , are uncorrelated).

Substituting equations (3) and (4) into equation (2), we obtain the reduced outcome model,

$$Y = \mu_Y^* + \beta X + \psi M + f^*(\mathbf{U}, \mathbf{V}, \mathbf{H}) + \varepsilon_Y^* \quad (5)$$

where $\beta = \eta\lambda$, $\psi = \phi\tau + \phi\delta\lambda$, $\mu_Y^* = \mu_Y + \lambda\mu_S + (\tau + \delta\lambda)\mu_P$, $\varepsilon_Y^* = \varepsilon_Y + \lambda\varepsilon_S + (\tau + \delta\lambda)\varepsilon_P$, and $f^*(\mathbf{U}, \mathbf{V}, \mathbf{H}) = f_Y(\mathbf{U}, \mathbf{V}, \mathbf{H}) + \lambda f_S(\mathbf{H}) + (\tau + \delta\lambda)f_P(\mathbf{H})$. Equation (5) represents the outcome model in Fig. 1b.

Because the instrumental variable Q is randomized, and hence statistically independent of any variables that are not directly or indirectly influenced by Q , it follows from equation (5) and standard properties of the covariance operator that,

$$\begin{aligned} \text{Cov}(Q, Y) &= \text{Cov}(Q, \mu_Y^*) + \beta \text{Cov}(Q, X) + \psi \text{Cov}(Q, M) \\ &\quad + \text{Cov}(Q, f^*(\mathbf{U}, \mathbf{V}, \mathbf{H})) + \text{Cov}(Q, \varepsilon_Y^*) \\ &= \psi \text{Cov}(Q, M), \end{aligned} \quad (6)$$

since $Q \perp\!\!\!\perp \mu_Y^*$, $Q \perp\!\!\!\perp X$, $Q \perp\!\!\!\perp f^*(\mathbf{U}, \mathbf{V}, \mathbf{H})$, and $Q \perp\!\!\!\perp \varepsilon_Y^*$, and the respective covariance terms are 0 (here, the symbol $\perp\!\!\!\perp$ stands for statistical independence). Therefore, ψ can be identified as,

$$\psi = \frac{\text{Cov}(Q, Y)}{\text{Cov}(Q, M)}, \quad (7)$$

as long as $\text{Cov}(Q, M) \neq 0$ (in practice, this condition is met if the psychological encouragement treatment can effectively manipulate the desire for improved symptoms, which, by its turn influences the emotional state, M). Note that the derivation of equation (7) required the key assumption that Y is linearly associated with M , and that the influence of residual errors and unmeasured confounders on Y is additive.

Now, if we let $R = Y - \psi M$ represent the residual of the outcome variable, after removal of the placebo effect, then we can rewrite equation (5) as,

$$R = \mu_Y^* + \beta X + f^*(\mathbf{U}, \mathbf{V}, \mathbf{H}) + \varepsilon_Y^*. \quad (8)$$

Because Z is also randomized, it follows from equation (8) and the properties of the covariance operator that,

$$\begin{aligned} \text{Cov}(Z, R) &= \text{Cov}(Z, \mu_Y^*) + \beta \text{Cov}(Z, X) + \text{Cov}(Z, f^*(\mathbf{U}, \mathbf{V}, \mathbf{H})) + \text{Cov}(Z, \varepsilon_Y^*) \\ &= \beta \text{Cov}(Z, X), \end{aligned} \quad (9)$$

since $Z \perp\!\!\!\perp \mu_Y^*$, $Z \perp\!\!\!\perp f^*(\mathbf{U}, \mathbf{V}, \mathbf{H})$, and $Z \perp\!\!\!\perp \varepsilon_Y^*$. Hence, the treatment effect β can be identified as,

$$\beta = \frac{\text{Cov}(Z, R)}{\text{Cov}(Z, X)}, \quad (10)$$

as long as $\text{Cov}(Z, X) \neq 0$ (in practice, this condition is met whenever there is some degree of compliance between assigned and received treatments). Note that in the derivation of equation (10) we assumed a linear relation between Y and M in the derivation of the residual R , as well as, that R and X (and, hence, Y and X) are also associated via a linear relationship.

Note that in addition to the three core assumptions required by an IV¹⁶ described above, the identification of the causal effects ψ and β require that: these effects are constant across the population; that the relationships between Y and M and between Y and X are linear; and that the contribution of the error terms, ε^* , and of the confounders, $f^*(\mathbf{U}, \mathbf{V}, \mathbf{H})$, to the response is additive (although it is not necessary to make any explicit assumptions about the form of the relationship between the confounders and the response).

Additionally, because the covariance operator only captures linear associations between two variables, and it is possible that two variables have zero covariance when the causal influence of the first variable on the second is mediated by a non-linear mechanism, the identification results in equations (7) and (10) require the additional assumptions that Q is linearly associated with M and Y , and that Z is linearly associated with X and R . We point out, however, that for binary instruments these additional assumptions can be relaxed since it can be shown (see the section on the non-parametric identification of causal effects for binary instruments) that the large sample estimators of the non-parametric average causal effects (derived without making the linearity assumption) are proportional to the respective sample covariance estimators, that is,

$$\widehat{\text{ACE}}(Q \rightarrow Y) = \widehat{E}(Y|do(Q = 1)) - \widehat{E}(Y|do(Q = 0)) \propto \widehat{\text{Cov}}(Q, Y), \quad (11)$$

$$\widehat{\text{ACE}}(Q \rightarrow M) = \widehat{E}(M|do(Q = 1)) - \widehat{E}(M|do(Q = 0)) \propto \widehat{\text{Cov}}(Q, M), \quad (12)$$

$$\widehat{\text{ACE}}(Z \rightarrow R) = \widehat{E}(R|do(Z = 1)) - \widehat{E}(R|do(Z = 0)) \propto \widehat{\text{Cov}}(Z, R), \quad (13)$$

$$\widehat{\text{ACE}}(Z \rightarrow X) = \widehat{E}(X|do(Z = 1)) - \widehat{E}(X|do(Z = 0)) \propto \widehat{\text{Cov}}(Z, X). \quad (14)$$

This observation shows that, at least for binary instruments (and when the sample size is large) it is not possible for two variables to have zero covariance when the causal influence of the first variable on the second is mediated by non-linear mechanisms.

Non-parametric identification of causal effects for binary instruments. In the following we derive a large sample non-parametric estimator of the causal effects of a binary instrumental variable using Pearl's interventional calculus⁴¹. But first we introduce some notation and provide a brief background.

Under the mechanism-based account of causation, the statistical information encoded in the joint probability distribution is supplemented with a causal DAG encoding the qualitative description of our assumptions about the causal relations between the variables. The joint probability distribution factorizes according to the causal graph,

$$P(x_1, \dots, x_p) = \prod_j P(x_j | pa(x_j)), \quad (15)$$

where each element, $P(x_j | pa(x_j))$, represents an autonomous mechanism describing the relationship between variable X_j and its parents. In this framework, causation means predicting the consequences of an intervention over a set of variables in the DAG, where intervention is expressed as a “surgery” on the equations and associated causal graph. We adopt the *do* operator notation to distinguish $P(y | do(X=x))$ from $P(y | X=x)$, where the former quantity describes the post-intervention distribution of Y given that the value of X was set to x by an external intervention, while the latter represents the usual conditional distribution of Y given that we observed the value of X to be equal to x (and is denoted the observational or pre-intervention distribution). For interventions over a single variable, the relationship between the pre-intervention and post-intervention distributions is given by the truncated factorization formula,

$$P(x_1, \dots, x_p | do(X_k = x'_k)) = \prod_{j \neq k} P(x_j | pa(x_j)) \mathbb{1}\{x_k = x'_k\}, \quad (16)$$

where the removal of the equation $P(x_k | pa(x_k))$ from the product in equation (16), and the replacement of x_k by x'_k in all elements $P(x_j | pa(x_j))$ for which X_k is a parent of X_j , formalizes what is meant by an “intervention surgery”. The average causal effect of a binary variable A on a variable B is defined as,

$$\text{ACE}(A \rightarrow B) = E(B | do(A = 1)) - E(B | do(A = 0)). \quad (17)$$

where the expectation is taken with respect to the post-intervention distribution $P(B | do(A = a))$. We say that the causal effect of A on B is identifiable if the post-intervention distribution $P(A | do(A = a))$ (and hence the $\text{ACE}(A \rightarrow B)$ quantity) is a function of observed variables only.

Now we show that the large sample non-parametric estimator of the causal effects of a binary instrumental variable on one of its descendent variables, is proportional to the respective sample covariance estimator. We illustrate the derivation using the average causal effect of Q on M , but the same exact rationale applies to the derivation of the causal effects of Q on Y , Z on X , and Z on R .

Let \mathcal{G} represent a DAG for which the core IV assumptions *i* to *iii* described above hold, but otherwise arbitrary. Note that, in this case, Q will always be an exogenous variable in \mathcal{G} (i.e., Q has no parents in \mathcal{G}). Let \mathbf{V} represent the set of all variables in \mathcal{G} , and $\mathbf{A} = \mathbf{V} \setminus \{M, Q\}$. Since Q is an exogenous variable in \mathcal{G} , we can factor the joint distribution of \mathbf{V} as,

$$P(M, \mathbf{A}, Q) = P(M, \mathbf{A} | Q = q)P(Q = q). \tag{18}$$

Although the conditional joint distribution, $P(M, \mathbf{A} | Q = q)$, can be further factorized according to \mathcal{G} , we don't need to specify the factorization explicitly when determining the post-intervention distribution for the intervention $do(Q = q')$, since application of the truncated factorization formula reduces to removing $P(Q = q)$, and replacing $Q = q$ by $Q = q'$ in the remaining conditional distributions, so that,

$$P(M, \mathbf{A} | do(Q = q')) = P(M, \mathbf{A} | Q = q'), \tag{19}$$

independent of how $P(M, \mathbf{A} | Q = q')$ can be further factorized. The marginal post-intervention distribution is given by,

$$P(M | do(Q = q')) = \sum_{\mathbf{A}} P(M, \mathbf{A} | Q = q') = P(M | Q = q'), \tag{20}$$

where the summation over \mathbf{A} is simply a notation for all the summations or integrations over each one of the variables in the set \mathbf{A} .

The average causal effect of Q on M is then given by,

$$\begin{aligned} ACE(Q \rightarrow M) &= E(M | do(Q = 1)) - E(M | do(Q = 0)) \\ &= E(M | Q = 1) - E(M | Q = 0), \end{aligned} \tag{21}$$

where the second equality follows from (20). A large sample non-parametric estimator of the expectation $E(M | Q = q')$ is given by,

$$\frac{\sum_{i=1}^n M_i \mathbb{I}\{Q_i = q'\}}{\sum_{i=1}^n \mathbb{I}\{Q_i = q'\}}, \tag{22}$$

so that,

$$\begin{aligned} \widehat{ACE}(Q \rightarrow M) &= \frac{\sum_{i=1}^n M_i \mathbb{I}\{Q_i = 1\}}{\sum_{i=1}^n \mathbb{I}\{Q_i = 1\}} - \frac{\sum_{i=1}^n M_i \mathbb{I}\{Q_i = 0\}}{\sum_{i=1}^n \mathbb{I}\{Q_i = 0\}} \\ &= \frac{\sum_{i=1}^n M_i Q_i}{\sum_{i=1}^n Q_i} - \frac{\sum_{i=1}^n Q_i (1 - Q_i)}{\sum_{i=1}^n (1 - Q_i)} \\ &= \frac{n^{-1} \sum_{i=1}^n Q_i M_i - (n^{-1} \sum_{i=1}^n Q_i)(n^{-1} \sum_{i=1}^n M_i)}{(n^{-1} \sum_{i=1}^n Q_i)(1 - n^{-1} \sum_{i=1}^n Q_i)} \propto \widehat{Cov}(Q, M). \end{aligned} \tag{23}$$

Two-step estimation procedure. Adopting a method of moments approach, we have that a consistent estimator for ψ is given by,

$$\hat{\psi}_{IV} = \frac{\widehat{Cov}(Q, Y)}{\widehat{Cov}(Q, M)} = \frac{\frac{1}{n} \sum_{k=1}^n Q_k Y_k - \left(\frac{1}{n} \sum_{k=1}^n Q_k\right) \left(\frac{1}{n} \sum_{k=1}^n Y_k\right)}{\frac{1}{n} \sum_{k=1}^n Q_k M_k - \left(\frac{1}{n} \sum_{k=1}^n Q_k\right) \left(\frac{1}{n} \sum_{k=1}^n M_k\right)}. \tag{24}$$

Note that the above placebo effect estimator requires measurements of M , but not of E or D . We point out, however, that if expectation and desire measurements are also available, then we can evaluate the validity of the desire-expectation model for the data at hand by checking whether the E , D , and I variables are able to predict the M measurements. We can also assess the effectiveness of the psychological treatment in influencing desire for better symptoms by estimating $Cor(Q, D)$.

Direct estimation of the treatment effect in equation (10) using an IV estimator is unfeasible, as it would involve the unmeasured quantities $R_k = Y - \psi M_k$. Therefore, in order to obtain a consistent estimator of the treatment effect, we adopt a two-step approach where we first estimate R_k as $\hat{R}_k = Y_k - \hat{\psi}_{IV} M_k$, for $k = 1, \dots, n$, and then estimate β using,

$$\hat{\beta}_{2sIV} = \frac{\widehat{Cov}(Z, \hat{R})}{\widehat{Cov}(Z, X)} = \frac{\frac{1}{n} \sum_{k=1}^n Z_k \hat{R}_k - \left(\frac{1}{n} \sum_{k=1}^n Z_k\right) \left(\frac{1}{n} \sum_{k=1}^n \hat{R}_k\right)}{\frac{1}{n} \sum_{k=1}^n Z_k X_k - \left(\frac{1}{n} \sum_{k=1}^n Z_k\right) \left(\frac{1}{n} \sum_{k=1}^n X_k\right)}. \tag{25}$$

Note that the IV estimators in equations (24) and (25) can produce highly inflated estimates when $\widehat{Cov}(Q, M) \approx 0$ and $\widehat{Cov}(Z, X) \approx 0$. Hence, in practice, it is important to check the assumptions that psychological encouragement influences the emotion levels, and that the compliance between assigned and received treatments is not negligible.

Randomization tests for $H_0: \psi = 0$ and $H_0: \beta = 0$. We implemented randomization tests⁴² for testing the presence of a placebo effect ($H_0: \psi = 0$ versus $H_1: \psi \neq 0$), and of a treatment effect ($H_0: \beta = 0$ versus $H_1: \beta \neq 0$). The randomization null distribution for the placebo effect is generated by evaluating the statistic $\hat{\psi}_{IV}$ in equation (24)

on a large number of shuffled versions of the data, where the Y_k measurements are shuffled relative to the (Q_k, M_k) measurements (whose connection is kept intact in order to preserve the association between the Q and M variables). The randomization null for treatment effect is generated by first calculating the residuals, $\hat{R}_k = Y_k - \hat{\psi}_{IV} M_k$, where $\hat{\psi}_{IV}$ is computed in the observed (not permuted) data, and then evaluating the statistic $\hat{\beta}_{2sIV}$ in equation (25) in shuffled data sets, where R_k is shuffled relative to (Z_k, X_k) data (whose connection is kept intact to preserve the association between Z and X). These randomization tests are non-parametric procedures and don't make any distributional assumptions about the data. However, because the identification of the causal effects assumes a linear relation between Y and X and M , the validity of the tests is still contingent on this assumption.

Randomization confidence intervals. Here we describe how to build confidence intervals for placebo and treatment effects using the p-values from randomization tests^{42,43}. Throughout this section we use θ to represent either the placebo effect, ψ , or the treatment effect, β . The procedure is straight-forward but requires a considerable amount of computation (which, nonetheless, can be easily parallelized). Assume for a moment that randomization tests for testing $H_0: \theta = \theta_j$ against one-sided alternative hypotheses $H_1: \theta < \theta_j$ and $H_1: \theta > \theta_j$ are available. Exploring the correspondence between confidence intervals and hypothesis tests, we obtain a $100(1 - 2\alpha)\%$ confidence interval (CI) for θ by searching for a lower bound value, θ_L , such that $H_0: \theta = \theta_L$ is rejected in favor of $H_1: \theta > \theta_L$ at a significance α , and by searching for an upper bound value, θ_U , such that $H_0: \theta = \theta_U$ is rejected in favor of $H_1: \theta < \theta_U$ at the same significance level⁴³.

While an efficient algorithm for finding CI bounds has been proposed⁴³, the approach requires the specification of the significant level before hand. In order to avoid this constraint, we generate a one-sided randomization p-value profile which can be used to determine the $100(1 - 2\alpha)\%$ CI for any desired α level. This p-value profile is generated as follows: (i) compute the observed placebo or treatment effect estimate, $\hat{\theta}$; (ii) for each $\theta_j < \hat{\theta}$, in a grid of decreasing θ_j values, compute the randomization p-value from the one-sided test $H_0: \theta = \theta_j$ vs $H_1: \theta > \theta_j$; (iii) repeat step ii until a p-value equal to zero is reached; (iv) for each $\theta_j > \hat{\theta}$, in a grid of increasing θ_j values, compute the p-value from the one-sided test $H_0: \theta = \theta_j$ vs $H_1: \theta < \theta_j$; (v) repeat step iv until a randomization p-value equal to zero is found.

Before we explain how to generate null distributions for placebo effects different from zero, consider first the intention-to-treat (ITT) estimator,

$$\widehat{ITT}_{\psi} = \frac{\sum_{k=1}^n Y_k \mathbb{1}\{Q_k = 1\}}{\sum_{k=1}^n \mathbb{1}\{Q_k = 1\}} - \frac{\sum_{k=1}^n Y_k \mathbb{1}\{Q_k = 0\}}{\sum_{k=1}^n \mathbb{1}\{Q_k = 0\}} = \frac{\widehat{\text{Cov}}(Q, Y)}{\widehat{\text{Var}}(Q)}. \quad (26)$$

Instead of directly generating a randomization distribution under the null $H_0: \psi = \psi_j$, we generate a randomization distribution under the equivalent null hypothesis that the intention-to-treat effect is equal to $\psi_j K_1$, where $K_1 = \widehat{\text{Cov}}(Q, M) / \widehat{\text{Var}}(Q)$ is constant across all permutations of the response data used in the construction of the randomization null. (Note that, because $\widehat{ITT}_{\psi} = K_1 \hat{\psi}$ the randomization tests based on $\hat{\psi}$ and \widehat{ITT}_{ψ} estimators produce exactly the same p-value if we use the same permutations of the response data in the construction of their null distributions).

The practical advantage of the test based on \widehat{ITT}_{ψ} effects is that it amounts to a simple two sample location problem for testing whether the difference in average response between the assigned treatment (psychological encouragement) and assigned control (no encouragement) groups is equal to $\psi_j K_1$. The implementation of randomization tests for this two sample location problem is straight-forward⁴³: we only need to add $\psi_j K_1$ for each Y_k data point in the assigned control group (i.e., k for which $Q_k = 0$), while leaving the response data from the assigned treatment group, $Q_k = 1$, unchanged, and then run a randomization test for testing the null hypothesis that the \widehat{ITT}_{ψ} effect is equal to zero, against the alternative one-sided hypothesis that it is positive, and against the alternative that it is negative.

Similarly, for the treatment effects we consider the two-step ITT estimator,

$$\widehat{ITT}_{\beta} = \frac{\sum_{k=1}^n \hat{R}_k \mathbb{1}\{Z_k = 1\}}{\sum_{k=1}^n \mathbb{1}\{Z_k = 1\}} - \frac{\sum_{k=1}^n \hat{R}_k \mathbb{1}\{Z_k = 0\}}{\sum_{k=1}^n \mathbb{1}\{Z_k = 0\}} = \frac{\widehat{\text{Cov}}(Z, \hat{R})}{\widehat{\text{Var}}(Z)}, \quad (27)$$

and generate randomization distributions under the equivalent null hypotheses $H_0: ITT_{\beta} = \beta_j K_2$, where $K_2 = \widehat{\text{Cov}}(Z, X) / \widehat{\text{Var}}(Z)$, by simply adding $\beta_j K_2$ for each \hat{R}_k data point in the assigned control group, $Z_k = 0$, while leaving the residual data from the assigned treatment group, $Z_k = 1$, unchanged (and then testing for the null that the \widehat{ITT}_{β} is equal to zero, against the alternative one-sided hypotheses that it is positive and the alternative that it is negative).

Adjustment for observed confounders. If measured confounders influencing both X and Y are available, it is possible to adjust for them by simply working with the residuals of X and Y (computed by separately regressing X and Y on the measured confounders). Similarly, if measured confounders influencing both M and Y are available, it is possible to adjust for them by working with the respective residuals. Note that even though, in theory, this type of adjustment is unnecessary, given that IVs allow for the consistent estimation of the causal effect in the presence of observed and unobserved confounders, it turns out that, in practice, it is possible to improve the power to detect causal effects by adjusting for observed confounders. Supplementary Fig. 6 shows an illustrative example, where the placebo effect estimator is adjusted by the treatment variable (which corresponds to a measured confounder of the placebo effect in unblinded trials).

Regression based estimators and tests. We compare the proposed IV estimators, and their respective randomization tests, to standard estimators and analytical hypothesis tests based on the linear regression of the outcome variable, Y , on both the received treatment, X , and emotion level, M , according to the model, $Y = \mu_Y + \beta X + \psi M + \varepsilon_Y$. Under this regression based approach, we estimate β and ψ using ordinary least squares,

and test the null hypotheses $H_0:\psi=0$ and $H_0:\beta=0$ using standard t-tests. In our simulations (described next), we generate data using gaussian errors, so that the distributional assumptions underlying the analytical t-tests are met.

Simulation experiments details. We simulated data from blinded and unblinded settings, in the presence or absence of confounding, according to the models presented in Fig. 2. For each of these settings, we run 4 separate simulation studies generating data: (i) under the null hypothesis that both treatment and placebo effect are zero, $H_0:\beta=0$ and $H_0:\psi=0$; (ii) under the alternative for treatment effects, $H_1:\beta \neq 0$, but null for placebo effects, $H_0:\psi=0$; (iii) the other way around, $H_0:\beta=0$ and $H_1:\psi \neq 0$; and (iv) under the alternative for both treatment and placebo effects, $H_1:\beta \neq 0$ and $H_1:\psi \neq 0$.

Each simulated data set was generated as follows. The IVs Z and Q were sampled from Bernoulli(1/2) distributions. All confounding variables were sampled from Normal(0, 1) distributions. The binary variables X , E , and D were generated by the threshold models,

$$X = \mathbb{1}\{\theta_{XZ}Z + \theta_{XU}U + \theta_{XC_1}C_1 + \theta_{XC_2}C_2 + \theta_{XC_3}C_3 + \varepsilon_X > 0\}, \quad (28)$$

$$E = \mathbb{1}\{\theta_{EX}X + \theta_{EC_1}C_1 + \theta_{EL_1}L_1 + \theta_{EV_2}V_2 + \theta_{EL_3}L_3 + \varepsilon_E > 0\}, \quad (29)$$

$$D = \mathbb{1}\{\theta_{DQ}Q + \theta_{DV_1}V_1 + \theta_{DC_2}C_2 + \theta_{DL_2}L_2 + \theta_{DL_3}L_3 + \varepsilon_D > 0\}, \quad (30)$$

where ε_X , ε_E , and ε_D were sampled from Normal(0, 1) distributions. The interaction I was generated as the product of E and D . Finally, the emotion and outcome data were generated from the linear models,

$$M = \theta_{ME}E + \theta_{MD}D + \theta_{MI}I + \theta_{ML_1}L_1 + \theta_{ML_2}L_2 + \theta_{MC_3}C_3 + \theta_{MV_3}V_3 + \varepsilon_M, \quad (31)$$

$$Y = \theta_{YX}X + \theta_{YM}M + \theta_{YU}U + \theta_{YV_1}V_1 + \theta_{YV_2}V_2 + \theta_{YV_3}V_3 + \varepsilon_Y, \quad (32)$$

where ε_M and ε_Y were sampled from Normal(0, 1) distributions. (Note that the explicit form of the desire-expectation model of emotions is unimportant, as the estimator for ψ depends on the observed values of M , but not of E , D , and I , and does not require an explicit description of the functional relationships between M and E , D , and I . Hence, for simplicity, we adopt a simple linear relation, even though more sophisticated relations could have been used).

Each simulation experiment comprised 10,000 distinct synthetic data sets. Each simulated data set was generated using a unique combination of simulation parameter values. In order to select parameter values spread as uniformly as possible over the entire parameter range we employed a Latin hypercube design⁴⁴, optimized according to the maximin distance criterium⁴⁵, in the determination of the parameter values used on each of the 10,000 simulated data sets for each simulation experiment.

We selected wide ranges for all model parameters. Explicitly, the parameters representing the effect of confounders on the observed variables (namely, θ_{XU} , θ_{XC_1} , θ_{XC_2} , θ_{XC_3} , θ_{EC_1} , θ_{EL_1} , θ_{EV_2} , θ_{EL_3} , θ_{DV_1} , θ_{DC_2} , θ_{DL_2} , θ_{DL_3} , θ_{ML_1} , θ_{ML_2} , θ_{MC_3} , θ_{MV_3} , θ_{YU} , θ_{YV_1} , θ_{YV_2} , and θ_{YV_3}) were selected in the range $[-2, 2]$ for the simulations under the influence of confounders, but were set to 0 in the simulations under unconfounded conditions. The effect of Z on X (θ_{XZ}), and of Q on D (θ_{DQ}), as well as, the effects of E , D , and I on M (θ_{ME} , θ_{MD} , and θ_{MI}) were selected in the range $[1, 2]$. The effect of X on E (θ_{EX}) was set to 0 in the blinded setting simulations, and selected in the range $[1, 2]$ in the unblinded simulations. The treatment effect (β) and the placebo effect (ψ) parameters were set to 0 in the simulations under the null hypothesis, and were selected in the range $[-2, 2]$ for the simulations under the alternative hypothesis. The range of sample size parameter, n , was set to realistic values we expect to see in practice, $\{100, 101, \dots, 1000\}$.

For any fixed significance level α , the empirical type I error rate was computed as the fraction of the p-values smaller than α over the data sets simulated under the null hypothesis, whereas the empirical power was calculated as the fraction of p-value smaller than α over data sets generated under the alternative hypothesis.

Code availability. The R code⁴⁶ implementing the IV approach, and used in the generation of the simulation results and figures, is available at: <https://www.synapse.org/DisentanglingTreatmentAndPlacebo>.

References

- Price, D. D., Finniss, D. G. & Benedetti, F. A comprehensive review of the placebo effect: recent advances and current thought. *Ann. Rev. Psychol.* **50**, 565–590 (2008).
- Finniss, D. G., Kaptchuk, T. J., Miller, F. & Benedetti, F. Biological, clinical, and ethical advances of placebo effects. *Lancet.* **375**, 686–695 (2010).
- Benedetti, F. Biological, clinical, and ethical advances of placebo effects. *Neuron.* **84**, 623–637 (2014).
- Fields, H. L. State-dependent opioid control of pain. *Nat. Rev. Neuro.* **5**, 565–575 (2004).
- Petrovic, P. *et al.* Placebo in emotional processing - induced expectations of anxiety relief activate a generalized modulatory network. *Neuron.* **46**, 957–969 (2005).
- Zubieta, J. K. *et al.* Regional μ -opioid receptor regulation of sensory and affective dimensions of pain. *Science.* **293**, 311–315 (2001).
- Price, D. D. & Barrell, J. J. Some general laws of human emotion: interrelationships between intensities of desire, expectation, and emotional feeling. *J. Personality.* **52**, 389–409 (1984).
- Price, D. D., Barrell, E. J. & Barrell, J. J. A quantitative-experiential analysis of human emotions. *Motivation and Emotion.* **9**, 19–38 (1985).
- Price, D. D. Assessing placebo effects without placebo groups: an untapped possibility? *Pain.* **90**, 201–203 (2001).
- Vase, L., Robison, M. E., Verne, G. N. & Price, D. D. The contributions of suggestion, expectancy and desire to placebo effect in irritable bowel syndrome patients. *Pain.* **105**, 17–25 (2003).
- Verne, G. N., Robison, M. E., Vase, L. & Price, D. D. Reversal of visceral and cutaneous hyperalgesia by local recatal anesthesia in irritable bowel syndrome (IBS) patients. *Pain.* **105**, 223–230 (2003).

12. Bausell, R. B., Lao, L., Bergman, S., Lee, W. L. & Berman, B. M. Is acupuncture analgesia and expectancy effect? Preliminary evidence based on participants' perceived assignments in two placebo controlled trials. *Eval. Health Professions*. **28**, 9–26 (2005).
13. McRae, C. *et al.* Effects of perceived treatment on quality of life and medical outcomes in a double-blind placebo surgery trial. *Arch. Gen. Psychiatry*. **61**, 412–420 (2004).
14. Angrist, J. & Krueger, A. Instrumental variables and the search for identification: from supply and demand to natural experiments. *J. Econ. Persp.* **15**, 69–85 (2001).
15. Bowden, R. J. & Turkington, D. A. *Instrumental Variables* (Cambridge University Press, 1990).
16. Didelez, V., Meng, S. & Sheehan, N. A. Assumptions of IV methods for observational epidemiology. *Stat. Science*. **25**, 22–40 (2010).
17. Levine, J. D., Gordon, N. C. & Fields, H. L. The mechanisms of placebo analgesia. *Lancet*. **2**, 654–657 (1978).
18. Benedetti, F. The opposite effects of the opiate antagonist naloxone and the cholecystokinin antagonist proglumide on placebo analgesia. *Pain*. **64**, 535–543 (1996).
19. Grevert, P., Albert, L. H. & Goldstein, A. Partial antagonism of placebo analgesia by naloxone. *Pain*. **16**, 129–143 (1983).
20. Levine, J. D. & Gordon, N. C. Influence of the method of drug administration on analgesic response. *Nature*. **312**, 755–756 (1984).
21. Amanzio, M., Pollo, A., Maggi, G. & Benedetti, F. Response variability to analgesics: a role for non-specific activation of endogenous opioids. *Pain*. **90**, 205–215 (2001).
22. Benedetti, F., Amanzio, M. & Maggi, G. Potentiation of placebo analgesia by proglumide. *Lancet*. **346**, 1231 (1995).
23. Zhang, Z., Kotz, R. M., Wang, C., Ruan, S. & Ho, M. A causal model for joint evaluation of placebo and treatment-specific effects in clinical trials. *Biometrics*. **69**, 318–327 (2013).
24. Jamshidian, F., Hubbard, A. E. & Jewell, N. P. Accounting for perception, placebo and unmasking effects in estimating treatment effects in randomised clinical trials. *Stat. Methods Med. Res.* **23**, 293–307 (2014).
25. Liu, W. *et al.* Joint estimation of treatment and placebo effects in clinical trials with longitudinal blinding assessments. *J. Am. Statist. Assoc.* (in press) (2016).
26. Roy, J. Randomized treatment-belief trials. *Contemporary Clinical Trials*. **33**, 172–177 (2012).
27. Balke, A. & Pearl, J. Bounds on treatment effects for studies with imperfect compliance. *J. Am. Statist. Assoc.* **92**, 1171–1176 (1997).
28. Hernan, M. A. & Robins, J. M. Instruments for causal inference: an epidemiologist's dream? *Epidemiology*. **17**, 360–372 (2006).
29. Swanson, S. A. *et al.* Bounding the per-protocol effect in randomized trials: an application to colorectal cancer screening *Trials*. **16**, 541 (2015).
30. Swanson, S. A. & Hernan, M. A. How to report instrumental variable analyses (suggestions welcome). *Epidemiology*. **24**, 370–374 (2013).
31. Rubin, D. B. Bayesian inference for causal effects: the role of randomization. *Annals of Statistics*. **6**, 34–58 (1978).
32. Imbens, G. W. & Rubin, D. B. *Causal inference for statistics, social and biomedical sciences: an introduction* (Cambridge University Press, 2015).
33. Angrist, J., Imbens, G. W. & Rubin, D. B. Identification of causal effects using instrumental variables. *J. Am. Statist. Assoc.* **91**, 444–455 (1996).
34. Hernan, M. A., Hernandez-Diaz, S. & Robins, J. M. A structural approach to selection bias. *Epidemiology*. **15**, 615–625 (2004).
35. Swanson, S. A., Robins, J. M., Miller, M. & Herman, M. A. Selecting on treatment: a pervasive form of bias in instrumental variable analysis. *Am. J. Epidemiol.* doi: 10.1093/aje/kwu284 (2015).
36. Jackson, J. W. & Swanson, S. A. Toward a clearer portrayal of confounding bias in instrumental variable applications. *Epidemiology*. **26**, 498–504 (2015).
37. Davies, N. M. An even clearer portrayal of bias in observational studies? *Epidemiology*. **26**, 505–508 (2015).
38. Chaibub Neto, E. *et al.* Towards personalized causal inference of medication response in mobile health: an instrumental variable approach for randomized trials with imperfect compliance. *arXiv:1604.01055v2* (2016).
39. Schork, N. J. Personalized medicine: time for one-person trials. *Nature*. **520**, 609–611 (2015).
40. Topol, E. *The Creative Destruction of Medicine* (Basic Books, 2012).
41. Pearl, J. *Causality: Models, Reasoning, and Inference* (Cambridge University Press, 2000).
42. Ernst, M. D. Permutation methods: a basis for exact inference. *Stat. Science*. **19**, 676–685 (2004).
43. Garthwaite, P. H. Confidence intervals from randomization tests. *Biometrics*. **52**, 1387–1393 (1996).
44. Santner, T. J., Williams, B. J. & Notz, W. I. *The Design and Analysis of Computer Experiments* (Springer Verlag, 2003).
45. Johnson, M., Moore, L. & Ylvisaker, D. Minimax and maximin distance designs. *J. Stat. Plan. Infer.* **26**, 131–148 (1990).
46. R Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/> (2014).

Acknowledgements

This work was funded by the Robert Wood Johnson Foundation.

Author Contributions

E.C.N. conceived the method, implemented and analyzed the simulation experiments, and wrote the manuscript.

Additional Information

Supplementary information accompanies this paper at <http://www.nature.com/srep>

Competing financial interests: The author declares no competing financial interests.

How to cite this article: Chaibub Neto, E. Using instrumental variables to disentangle treatment and placebo effects in blinded and unblinded randomized clinical trials influenced by unmeasured confounders. *Sci. Rep.* **6**, 37154; doi: 10.1038/srep37154 (2016).

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



This work is licensed under a Creative Commons Attribution 4.0 International License. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>

© The Author(s) 2016