

**Paired  $2 \times 2$  Factorial Design for Treatment Effect  
Identification and Estimation in the Presence of Paired  
Interference and Noncompliance**

**Máté Kormos**

## Abstract

I propose an experimental design, the Paired  $2 \times 2$  Factorial Design, together with two strategies for treatment effect identification and estimation from a large sample of pairs, comprised of distinguishable members (e.g. couples with a healthy and an ill member), when (i) there is interference: the potential outcomes of a pair member do not only depend on the member's own treatment participation but on his/her partner's participation too; (ii) there is (endogenous) noncompliance: the units may not comply perfectly with their treatment encouragement; (iii) the experimenter can have two different binary treatments and two different outcomes of interest for the different members in a pair; (iv) units within a pair are allowed to coordinate their treatment participation based on their encouragement. The latter, number (iv), is my main contribution, as this has not been addressed by previous studies.

The first strategy uses only half of the sample to identify the effects of both treatments on both pair members, in certain complier subpopulations, under the usual instrumental variables assumptions. The second strategy, the main theoretical result, uses the whole sample to identify the same treatment effects, relying on the same set of assumptions.

The use of the design is illustrated, in theory, on a sample of married couples where one member suffers from depression. The treatment is an antidepressant for the depressed member and an educational program for the healthy partner.

# Contents

<b>List of Tables</b>	<b>v</b>
<b>List of Figures</b>	<b>vi</b>
<b>List of Theorems</b>	<b>vii</b>
<b>List of Lemmas</b>	<b>viii</b>
<b>List of Algorithms</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Application Domains, Examples . . . . .	2
<b>2 Previous Work</b>	<b>3</b>
<b>3 Identification</b>	<b>6</b>
3.1 Framework & Notation . . . . .	6
3.1.1 Statistical Viewpoint & Sampling . . . . .	6
3.1.2 Partial Interference . . . . .	7
3.1.3 Distinguishable Members . . . . .	7
3.1.4 Treatment Encouragement . . . . .	7
3.1.5 Treatment Participation . . . . .	8
3.1.6 Potential Outcomes . . . . .	9
3.2 Variables in the Sample . . . . .	10
3.3 Paired $2 \times 2$ Factorial Design . . . . .	10
3.4 Identifying Assumptions . . . . .	12
3.4.1 Exclusion . . . . .	13
3.4.2 Random Assignment . . . . .	13
3.4.3 i.i.d. Assignment . . . . .	13
3.4.4 One-sided Noncompliance . . . . .	13
3.4.5 Monotonicity . . . . .	13
3.4.6 Invertibility . . . . .	14
3.5 Half-sample Identification Strategy . . . . .	14
3.6 Full-sample Identification Strategy . . . . .	17
<b>4 Conclusion, Limitations and Further Research</b>	<b>18</b>

<b>References</b>	<b>19</b>
<b>Appendix A Distribution of Treatment Participation</b>	<b>21</b>
A.1 Conditions for Identifying Assumptions . . . . .	21
A.2 Implications of Identifying Assumptions . . . . .	22
<b>Appendix B Proof of Theorem 1</b>	<b>24</b>
<b>Appendix C Proof of Theorem 2</b>	<b>25</b>
C.1 Generic Notation . . . . .	25
C.2 Statement of Theorem 2 . . . . .	26
C.3 Computing $\mathbf{M}$ . . . . .	27
C.4 Computing $\mathbf{v}$ . . . . .	34
C.5 Computing $\mathbf{M}^{-1}\mathbf{v}$ . . . . .	42
C.6 Overview of the Proof . . . . .	52

## List of Tables

Table 1	Application examples . . . . .	3
Table 2	Output from <code>symbolic_inversion.py</code> . . . . .	43

## List of Figures

Figure 1	Half- and Full-sample strategies, with $(Z_A, Z_B)$ values in the nodes . . . . .	14
----------	---	----

## List of Theorems

<b>Theorem 1</b>	<b>Half-sample Identification</b> . . . . .	15
<b>Theorem 2</b>	<b>Full-sample Identification</b> . . . . .	17

## List of Lemmas

Lemma 1	Conditions for $A5$ Monotonicity . . . . .	21
Lemma 2	Conditions for $A6$ Invertibility . . . . .	22
Lemma 3	Implications of $A5$ Monotonicity . . . . .	22
Lemma 4	Implications of $A6$ Invertibility . . . . .	23
Lemma 5	Determinant of $M$ . . . . .	32
Lemma 6	Existence of $M^{-1}$ . . . . .	32
Lemma 7	Sufficient Condition for the Existence of $M^{-1}$ . . . . .	33



## List of Algorithms

Algorithm 1	Paired $2 \times 2$ Factorial Design . . . . .	11
Algorithm 2	Paired $2 \times 2$ Factorial Design example . . . . .	11

## 1. Introduction

*Suppose that we have a random sample of married couples where one member of each couple suffers from depression while the other member, the depressed's spouse, does not. Of interest are the effects of two binary treatments: an antidepressant for the depressed member, and an educational program about depression for the healthy spouse.*

In this example, a depressed person's outcome (say, symptoms of depression) might be improved alone by his/her partner's participation in the educational program even if the depressed member does not take the antidepressant. The same is true for the healthy partner, who might feel better due to his/her depressed spouse taking the antidepressant even if he/she does not attend the educational program. If so, the Stable Unit Treatment Value Assumption (SUTVA (Rubin, 1980)), made by most works on treatment effects in the framework of the Rubin Causal Model (Rubin, 1974), is violated as interference arises within predefined pairs of units, the married couples.<sup>1</sup>

Hence, if a researcher wants to make causal inference about the treatment effects of the antidepressant and the educational program in a randomised controlled trial (RCT), a suitable experimental design and methods addressing interference are needed. Furthermore, it is necessary for such methods to correct for the cases when the units do not comply with their treatment encouragement instructions because their expected improvement from doing so is low or negative (in short, endogenous noncompliance, or, somewhat imprecisely, imperfect compliance). For example, if the expected improvement from participating in the educational program is low for both members, then the healthy member is less likely to participate in the educational program when he/she is encouraged to do so.

Methods capable of handling interference and imperfect compliance in RCTs have just been gaining ground recently but are still rare. Moreover, to the best of my knowledge, there exists no such method which allows for within-pair discussion and coordination of treatment participation even though this is desirable in real-life applications. To fill this gap, I propose an easy-to-implement experimental design, built on a coin-tossing protocol, the Paired  $2 \times 2$  Factorial Design, and two associated large-sample identification and estimation strategies, which

---

<sup>1</sup>SUTVA states that the potential outcomes of a unit (the unit's outcome were it (not) treated) depends only on the unit's own participation in the treatment but not on others'.

- (i) take into account interference present only within pairs of units, where the members of a pair are distinguishable (e.g. healthy-ill, male-female, sister-brother members in each pair in the population); and
- (ii) enable the experimenter to have two, potentially member-specific binary treatments and two, potentially member-specific outcomes of interest (or one single outcome belonging to the pair); and
- (iii) allow for imperfect compliance, potentially of endogeneous nature, of the units with the treatment encouragement given by the experimenter; and
- (iv) allow for within-pair discussion and coordination of treatment participation.

The first strategy uses only half of the sample to identify the average treatment effects of both treatments on both pair members, in certain complier subpopulations, under the usual instrumental variables assumptions (Imbens and Angrist, 1994). The second strategy, the main theoretical result, uses the whole sample to identify the same treatment effects, relying on the very same assumptions. Unfortunately, neither of the strategies is appropriate for identifying the joint effect of the two treatments; only a spillover-like effect can be captured by the second one.

Regarding application, the Paired  $2 \times 2$  Factorial Design together with the two identification strategies constitute a ready-made tool for practical use. Implemented in Python, each step is feasible.<sup>2</sup> In the next section, I give some examples for cases when the design could be used.

### 1.1. Application Domains, Examples

The proposed design allows for the evaluation of two, potentially different, binary treatments on two, potentially different, outcomes belonging to distinguishable pair members (for notations and definitions see Section 3.1). Thus every case in which we have a population of pairs with *distinguishable* members is a valid domain for the application. This includes cases when the treatments are the same for both pair members and/or when the two outcomes of interest are the same for both pair members. Moreover, cases where there is only a single outcome belonging to a pair, i.e. *not* two identical outcomes belonging to the two members, are included as well.

In general there exists a variable along which we can distinguish between members in a pair. To mention probably the most prevalent one, there is age: in each

---

<sup>2</sup>Find the codes at [GitHub](#).

pair we are almost sure to be able to distinguish younger from older members. It is another question, however, that what sense it makes to do so, whether the interpretation is meaningful in this way. This question is to be answered by the experimenter based on the research topic.

Nevertheless, I list some possible examples in Table 1, where the distinction is natural and the design could be useful. Take, e.g. a policy evaluation issue regarding parental leave. We can distinguish between the mother and the father in each couple, and the treatment is the maternity and paternity leave, with new conditions. It is quite reasonable to assume that the parents discuss whether they want to stay at home with the baby given the chance. In this case, the proposed design could be useful to identify the effect of parental leave on the labour supply of both parents afterwards, which could affect some labour policy. Another research on parental leave could assess its effect on the social/psychological development of their child. This is an example for the case when the outcome of interest (development of children) is a single variable and we have one measurement of it per pair, not two, belonging to two parents. A typical economic example is principal-agent contracting. The treatments could be new principal and agent specific clauses in the contract, with outcomes being the satisfaction of the principal with the agent’s work and the performance, or stress-level, of the agent.

**Table 1: Application examples**

Population		Treatment		Outcome	
Member-A	Member-B	Treatment-A	Treatment-B	Outcome-A	Outcome-B
impoverished sister	impoverished brother	special educational program	special educational program	grades	grades
husband	wife	marriage therapy for husbands	marriage therapy for wives	marriage outcome	
mother	father	maternity leave	paternity leave	labour supply	labour supply
mother	father	maternity leave	paternity leave	child development	
sister	brother	vaccination	vaccination	health status	health status
principal	agent	contract	contract	satisfaction	performance

The remainder of this paper is organised as follows. Section 2 contains the novel components in relation to the existing literature; Section 3 describes the experimental design and the two identification strategies. Finally, I conclude in Section 4.

## 2. Previous Work

The idea of interference is not new, with many papers written on the subject based on either randomised trials or observational data (for a review see [VanderWeele et al.](#)

(2014)). Without any specific structure imposed on the nature of interference, [Halloran and Struchiner \(1991\)](#) conceptualise treatment effects of interest, estimands such as direct and indirect effects. [Halloran and Struchiner \(1995\)](#) refine these concepts aligning them with Rubin’s paradigm. [Rosenbaum \(2007\)](#) extends Fisher’s sharp null hypothesis of no effect for the interference setting.

[Sobel \(2006\)](#) introduces the notion of partial interference, on which my design is built so that interference may only occur within groups, pairs, in the population but not across them. Based on [Sobel \(2006\)](#), [Hudgens and Halloran \(2008\)](#) propose a two-stage randomised design to identify and estimate total, direct and indirect treatment effects unbiasedly and consistently and derive variances of the estimators.

As for observational studies, [Tchetgen and VanderWeele \(2012\)](#) contributes to the literature by examining finite sample inference with inverse probability weighting estimators, which are implemented in R by [Saul and Hudgens \(2017\)](#). While they do not take into account the within-group correlation structure of treatment participation, [Barkley et al. \(2017\)](#) do so with a group-specific random effect logit model.

In all the randomised trial approaches above, it is assumed that the units perfectly comply with their treatment assignment instruction, which may be unrealistic. Under SUTVA, in the no interference case, the instrumental variables (IV) framework of [Imbens and Angrist \(1994\)](#) provides a means to correct for this and makes it possible to identify treatment effects for the subpopulation who comply with the assignment. However, [Sobel \(2006\)](#) proves that these IV methods do not work under interference. Another, simulation-based, evidence for the undesirable IV properties for the  $2 \times 2$  factorial design is given by [Merrill and McClure \(2015\)](#), which however assumes that there is no interference effect of the treatments on the outcome.

A solution to this is presented by [Kang and Imbens \(2016\)](#) who develop a two-stage randomised design called Peer Encouragement Design. An important step to account for imperfect compliance, they put forward the idea of personalised treatment. Personalised treatment means that a unit’s participation can only depend on the unit’s own treatment encouragement, but not on other units’ treatment encouragement. As the authors point out, this might be violated if the units are able to discuss and coordinate their participation knowing about each other’s encouragement.

My method is vary similar to that of [Blackwell \(2017\)](#) who corrects for non-compliance in the  $2 \times 2$  factorial design with IV, when there are two treatments and imperfect compliance can occur along both treatments. However, in his framework

there is no paired structure, the level of analysis is the individual/person. In other words, interference arises from the interaction of the two treatments not from another individual's treatment (he makes the assumption of no interference across people). While my paired experiment setup can be thought of as having two treatments for the same person (we consider education of the partner as the second treatment for the depressed spouse) with imperfect compliance occurring along both treatments, Blackwell does not study a paired structure. Neglecting paired structure is not necessarily a problem when we are only interested in the outcome of either the depressed or the non-depressed person (or in a single outcome per pair), yet it is restrictive when both are of interest.

More importantly, Blackwell makes the same personalised treatment assumption as [Kang and Imbens \(2016\)](#), corresponding to the individual-level analysis, so that participation in 'treatment-A' only depends on the encouragement of 'treatment-A' but not on that of 'treatment-B'. In some cases this may be plausible but in the paired structure it rules out within-pair discussion as [Kang and Imbens \(2016\)](#) writes, so it is less credible.

As opposed to this, my proposed design and strategies relax this assumption by allowing for such discussion and coordination within a pair. To the best of my knowledge, there exists no other work permitting this yet, thus my work can be considered as a first step in this direction.

Compared to previous works on interference, I focus my attention on pairs, which is restrictive on one hand, but the results are more easily interpretable. This is achieved by allowing for member-specific treatments (antidepressant for the depressed member, educational program for the spouse) and also for different outcome of interests for the different members of the pairs (depression symptoms for the depressed member, some well-being/anxiety measure for the spouse). Previous studies typically define group-level averages to identify effects more easily; however this would not work with member-specific treatments and outcomes of interest.

### 3. Identification

In this section, after the general framework and notation is introduced (Section 3.1), I present the protocol of the Paired  $2 \times 2$  Factorial Design (Section 3.3), with which are associated two identification strategies. Following the discussion of the identifying assumptions (Section 3.4), the two strategies are described.

The first strategy uses only half of the sample to infer treatment effects (Half-sample Identification Strategy, see Section 3.5). This strategy is not a theoretical advancement per se, it is merely a ‘trick’ to see the problem as one in the traditional [Imbens and Angrist \(1994\)](#) IV-setting. Consequently, the advantage of it is that the assumptions of [Imbens and Angrist \(1994\)](#) applied to the interference setup are sufficient for identification (see Theorem 1). The drawback is that it throws away the information present in the other half of the sample.

The advantage of the second strategy, the main theoretical result (Full-sample Identification Strategy, see Section 3.6 and Theorem 2), is that it uses the whole sample to identify treatment effects, yet it does not require additional assumptions.

#### 3.1. Framework & Notation

##### 3.1.1 Statistical Viewpoint & Sampling

In analysing relationships between variables, multiple approaches can be taken depending on how we view our available data at hand. We can consider our available data as the whole population or a random sample drawn from the population.<sup>3</sup> Throughout the thesis, I take the view that there is an infinite population of pairs from which the available data are a single random sample drawn with replacement. The reason is that, even if the whole population is available, we are likely to want to sample from it due to cost considerations. Sampling randomly with replacement from the population is important because this way the sample is independently and identically distributed which renders the analysis simpler. Given this viewpoint, there are two possible scenarios.

The first scenario is the cleanest, statistically speaking. Suppose that we have the complete list of the pairs in the population of interest. *In principle*, we assign treatment to every single unit in the population in the same way as treatment assignment protocol is given in Algorithm 1 with the difference that this time we do this for the whole population, not only for the sample. Then we draw a random

---

<sup>3</sup>For further details on the implications of this on the analysis of randomised experiments see [Athey and Imbens \(2016\)](#)

sample of pairs with replacement (e.g. we independently generate random integers uniformly corresponding to the unique identifier of the pairs), and now we tell the members in the sampled pairs that what their in-principle treatment encouragement was. In this way, there is a theoretical chance that a pair shows up in the sample multiple times (because we sample with replacement); however the probability of this is near-zero.

In the second scenario, suppose that we do not have the list of the population, but we have a large number of available pairs. Now if we assume that this available set is itself a random sample drawn with replacement from the population of interest, we can apply the very same encouragement protocol as that in Algorithm 1 for this sample. In this way, however, the probability that a pair is present in the sample multiple times is exactly zero.

Either way, we end up with an independently and identically distributed large sample of size  $n$  from the infinite population of pairs for the analysis. Hence, every random variable introduced below is a single draw from their respective distribution and is indexed with  $i$ :  $i = 1, \dots, n$ . In general, I omit the index  $i$ , except for the introduction below and when it is necessary to write it out.

### 3.1.2 Partial Interference

According to partial interference, the interference may occur only within specific groups, pairs in this case, but not across them. Interference stands for possible interactions between potential outcomes and treatment participation (see Section 3.1.5) within a pair.

### 3.1.3 Distinguishable Members

The pair members are distinguishable; thus we have member-A and member-B in each pair of the population. Equivalently, there exists at least one variable along which we can deterministically distinguish member-A from member-B in the whole population of interest. In our example, the discriminating variable is whether someone is depressed (say, member-A) or not (member-B). Other examples for division: healthy-ill, doctor-patient, sister-brother, man-woman, younger-older pairs etc..

### 3.1.4 Treatment Encouragement

Let  $Z_{A_i} \in \{0, 1\}$  and  $Z_{B_i} \in \{0, 1\}$ ,  $i = 1, \dots, n$  be the binary random variables indicating the treatment encouragement (assignment) of member-A and member-B



in the  $i$ th pair of the  $n$ -large sample. Then  $Z_{Ai} = 1$  if and only if member-A in the  $i$ th pair is encouraged to take the member-A specific treatment, treatment-A; and  $Z_{Bi} = 1$  if and only if member-B in the  $i$ th pair is encouraged to take the member-B specific treatment, treatment-B. For example, if member-A is the depressed person,  $Z_{Ai} = 1$  and  $Z_{Bi} = 1$  means that in pair  $i$  we encourage the depressed person to take the antidepressant and his/her partner to enroll to the educational program.

### 3.1.5 Treatment Participation

Denoting the observable, actual treatment participation with  $D$ , we can distinguish between four binary random variables belonging to member-A and another four belonging to member-B. As a function of the treatment encouragements in pair  $i$ , this is compactly written as  $D_{Ai}(Z_{Ai}, Z_{Bi}) \in \{0, 1\}$  and  $D_{Bi}(Z_{Ai}, Z_{Bi}) \in \{0, 1\}$ . That is, we have  $D_{Ai}(00), D_{Ai}(10), D_{Ai}(01), D_{Ai}(11)$  and  $D_{Bi}(00), D_{Bi}(10), D_{Bi}(01), D_{Bi}(11)$ .<sup>4</sup> To ease notation, no comma is used to separate arguments when actual numbers (0,1) are used. It is important to pay attention to the order of the arguments: A comes first so that  $D_{Ai}(01)$  and  $D_{Bi}(01)$  describes the actual treatment participation of members when member-A is not encouraged to take his/her treatment ( $Z_{Ai} = 0$ ) and member-B is encouraged to take his/her treatment ( $Z_{Bi} = 1$ ).

Among the four variables, there is only one which is observable as there is only one treatment encouragement which can be given to pair  $i$ . In our example, we obviously cannot instruct the very same depressed person to take and not to take the antidepressant at the same time. Conveniently, the observable one among the four random variables can be written as

$$D_{Ai} = Z_{Ai}Z_{Bi}D_{Ai}(11) + Z_{Ai}(1 - Z_{Bi})D_{Ai}(10) + (1 - Z_{Ai})Z_{Bi}D_{Ai}(01) + (1 - Z_{Ai})(1 - Z_{Bi})D_{Ai}(00) \quad (1)$$

$$D_{Bi} = Z_{Ai}Z_{Bi}D_{Bi}(11) + Z_{Ai}(1 - Z_{Bi})D_{Bi}(10) + (1 - Z_{Ai})Z_{Bi}D_{Bi}(01) + (1 - Z_{Ai})(1 - Z_{Bi})D_{Bi}(00), \quad (2)$$

where the binary treatment encouragement  $(Z_{Ai}, Z_{Bi})$  ‘activates’ the appropriate random variables  $(D_{Ai}(Z_{Ai}, Z_{Bi}), D_{Bi}(Z_{Ai}, Z_{Bi}))$  corresponding to the given treatment encouragement.

---

<sup>4</sup>The personalised treatment of [Kang and Imbens \(2016\)](#) formally states that  $D_{Ai}(Z_{Ai}, Z_{Bi}) = D_{Ai}(Z_{Ai})$  and  $D_{Bi}(Z_{Ai}, Z_{Bi}) = D_{Bi}(Z_{Bi})$ .

Once again, suppose that member-A is the depressed person. Then by not giving the depressed person the antidepressant ( $Z_{Ai} = 0$ ) while encouraging his/her spouse to enroll to the educational program ( $Z_{Bi} = 1$ ), we observe  $D_{Ai}(01)$  and  $D_{Bi}(01)$  telling us whether the depressed member takes the pill and whether the non-depressed one enrolls to the program given the encouragement (no pill, education program).

### 3.1.6 Potential Outcomes

There are four potential outcomes in pair  $i$  for member-A:  $Y_{Ai}(D_{Ai}, D_{Bi}) \in \mathbb{R}^1$  and four for member-B:  $Y_{Bi}(D_{Ai}, D_{Bi}) \in \mathbb{R}^1$ . These are four-four random variables which characterise the outcome of interest as a function of the *actual treatment participation* of both members in the pair. They can be thought of as random variables describing the outcome in four states of the world, of which we can only witness one. Similarly to the treatment encouragement, the observable one of the four is written as

$$Y_{Ai} = D_{Ai}D_{Bi}Y_{Ai}(11) + D_{Ai}(1 - D_{Bi})Y_{Ai}(10) + (1 - D_{Ai})D_{Bi}Y_{Ai}(01) + (1 - D_{Ai})(1 - D_{Bi})Y_{Ai}(00) \quad (3)$$

$$Y_{Bi} = D_{Ai}D_{Bi}Y_{Bi}(11) + D_{Ai}(1 - D_{Bi})Y_{Bi}(10) + (1 - D_{Ai})D_{Bi}Y_{Bi}(01) + (1 - D_{Ai})(1 - D_{Bi})Y_{Bi}(00), \quad (4)$$

where the binary actual treatment participations  $(D_{Ai}, D_{Bi})$  select the observable one. It might be the case that  $Y_{Ai}(D_{Ai}, D_{Bi})$  and  $Y_{Bi}(D_{Ai}, D_{Bi})$  measures the same thing for the members, i.e. GPA of member-A and GPA of member-B. It might also be the case that there is only a single potential outcome per pair (e.g. marriage outcome of married couples; see Section 1.1). If so,  $Y_{Ai}(D_{Ai}, D_{Bi}) = Y_{Bi}(D_{Ai}, D_{Bi})$ , thus A-B indexing is unnecessary and the effects in Section 3.5 and 3.6 is to be interpreted accordingly.

In our example (member-A is the depressed),  $Y_{Ai}(10)$  and  $Y_{Bi}(10)$  are the outcomes of the pair members when the depressed actually takes the antidepressant and his/her spouse is not enrolled to the educational program;  $Y_{Ai}(11)$  and  $Y_{Bi}(11)$  are the outcomes when the depressed takes the pill and the non-depressed is enrolled to the program and so on.

### 3.2. Variables in the Sample

To understand how the proposed estimator works, a few more notations have to be introduced. Let  $y_{Ai}$  ( $y_{Bi}$ ) denote the sample analogue of the observable outcome for member-A (member-B); let  $\mathbf{d}_{Ai} \equiv [1, d_{Ai}, d_{Bi}, d_{Ai}d_{Bi}]' \in \{0, 1\}^{4 \times 1}$  and  $\mathbf{d}_{Bi} \equiv [1, d_{Bi}, d_{Ai}, d_{Ai}d_{Bi}]' \in \{0, 1\}^{4 \times 1}$ , where  $d_{Ai}$  ( $d_{Bi}$ ) is the sample-realisation of treatment participation of member-A (member-B) in pair  $i$ ; and last let  $\mathbf{z}_{Ai} \equiv [1, z_{Ai}, z_{Bi}, z_{Ai}z_{Bi}]' \in \{0, 1\}^{4 \times 1}$  and  $\mathbf{z}_{Bi} \equiv [1, z_{Bi}, z_{Ai}, z_{Ai}z_{Bi}]' \in \{0, 1\}^{4 \times 1}$ , where  $z_{Ai}$  ( $z_{Bi}$ ) is the sample-realisation of the treatment encouragement of member-A (member-B) in pair  $i$ . Broadcasting them to  $\mathbf{y}_A \equiv [y_{A1}, \dots, y_{An}]' \in \mathbb{R}^{n \times 1}$  and  $\mathbf{y}_B \equiv [y_{B1}, \dots, y_{Bn}]' \in \mathbb{R}^{n \times 1}$ ,  $\mathbf{D}'_A \equiv [\mathbf{d}_{A1}, \dots, \mathbf{d}_{An}] \in \{0, 1\}^{4 \times n}$  and  $\mathbf{D}'_B \equiv [\mathbf{d}_{B1}, \dots, \mathbf{d}_{Bn}] \in \{0, 1\}^{4 \times n}$ ,  $\mathbf{Z}'_A \equiv [\mathbf{z}_{A1}, \dots, \mathbf{z}_{An}] \in \{0, 1\}^{4 \times n}$  and  $\mathbf{Z}'_B \equiv [\mathbf{z}_{B1}, \dots, \mathbf{z}_{Bn}] \in \{0, 1\}^{4 \times n}$  facilitates more compact notation, with “ $'$ ” indicating the transpose.

### 3.3. Paired $2 \times 2$ Factorial Design

After having obtained a large random sample of pairs, the next step in the experiment is to assign treatments to the units (e.g. which depressed person is encouraged to take the antidepressant and which healthy person is encouraged to enroll to the educational program). The exact, algorithmic procedure of doing so is the experimental design/protocol, which is presented in this section.

In our case, a good starting point is the  $2 \times 2$  factorial design, which is suitable for exploring the interaction effect between two binary treatments on a single unit (Cheng, 2013). The depression example can be thought of as having two treatment for the same unit/person (we consider education of the partner as the second treatment for the depressed spouse) which gives rise to the factorial design.

The treatment encouragement protocol is specified in Algorithm 1, and is illustrated with the depression example in Algorithm 2. The encouragement mechanism is fairly simple and the advantage of it is that upon treatment assignment, the experimenter does not have to take into account the paired structure: it is only later, in the treatment effect estimation phase, when it is necessary to keep track of the paired structure. In contrast to the personalised treatment assumption of Kang and Imbens (2016), the experimenter does not have to care about whether the members in a pair know about each other’s encouragement or not, as discussion and coordination between them is allowed for in the estimation phase.

---

**Algorithm 1 Paired  $2 \times 2$  Factorial Design**

---

```
1: draw an  $n$ -large i.i.d. sample from the population of (member-A, member-B) pairs
2: for  $k=1:2n$  do ▷ for each unit,  $k$ , in the sample
3:   draw  $x$  from Bernoulli( $P$ ) i.i.d. with  $P = 0.5$ 
4:   if  $x \geq 0.5$  then
5:     if person is member-A then
6:       encouragement(person)  $\leftarrow$  ‘take treatment-A!’
7:     else
8:       encouragement(person)  $\leftarrow$  ‘take treatment-B!’
9:      $z_k \leftarrow 1$  ▷ dummy indicating treatment
10:  else
11:    encouragement(person)  $\leftarrow$  excluded from treatment
12:     $z_k \leftarrow 0$ 
```

**Output:**  $n$  pairs with member-specific treatment assignment:  $\{(z_{Ai}, z_{Bi})\}_{i=1}^n$

---

---

**Algorithm 2 Paired  $2 \times 2$  Factorial Design example**

---

```
1: draw an i.i.d. sample from the population of (depressed, not depressed) pairs
2: for each person in the sample do
3:   flip a fair coin:  $\mathbb{P}(\text{head}) \equiv P = 0.5$ 
4:   if head then
5:     if person is depressed then
6:       encouragement(person)  $\leftarrow$  ‘take the pill!’
7:     else
8:       encouragement(person)  $\leftarrow$  ‘enroll to educational program!’
9:   else
10:    encouragement(person)  $\leftarrow$  cannot access to pill/education
```

**Output:**  $n$  pairs with member-specific treatment assignment

---

The encouragement protocol leads to four groups of pairs along treatment assignment: based on the observed values of  $(Z_A, Z_B)$  we have  $\mathcal{G}_{Z_A, Z_B} : \mathcal{G}_{00}, \mathcal{G}_{10}, \mathcal{G}_{01}, \mathcal{G}_{11}$ . In our example (member-A is depressed), these are the pairs with treatment encouragement: (no pill, no education), (pill, education), (no pill, education), (pill, education), respectively. It is easy to see that for  $P = 0.5$  the groups consist of approximately the same number of pairs, i.e. the proportion of pairs in each group is roughly 25%, even though we neglected the paired structure during encouragement.<sup>5</sup> These groups play a role in the different identification strategies in Section 3.5 and 3.6.

---

<sup>5</sup>See the Monte Carlo simulation of the assignment in `assignment_mechanism.py`, which verifies the 25% proportion.

### 3.4. Identifying Assumptions

Once the treatments are assigned to the units based on the protocol in Algorithm 1, the experimenter observes whether they comply and participate in the treatment, and records their outcome of interest (e.g. symptoms of depression). Before analysing the data to identify treatment effects with the strategies in Section 3.5 or 3.6, it is crucial that the experimenter is aware of the underlying assumptions. In this section, the Identifying Assumptions, i.e. the sufficient conditions to identify treatment effects, are outlined. The two (half- and full-sample) identification strategies depend on the same set of assumptions:  $A1 - A6$ . First, I describe the assumptions formally, and then in a more intuitive way.

Some assumptions address the treatment encouragement and thus are completely under the control of the experimenter; they are met by following Algorithm 1 (these assumptions are not in *italic*). Other assumptions concern the participation behaviour of the subjects and hence are not under the control of the experimenter (indicated in *italic*). In between these is *One-sided noncompliance*, which depends on both the experimenter and the nature of the treatments. An in-depth view on assumptions concerning the distribution of treatment participation ( $A4 - A8$ ) is provided in Appendix A.

**Identifying Assumptions** Usual IV assumptions, extended for interference:

**A1** Exclusion: 
$$Y_A(D_A, D_B, Z_A, Z_B) = Y_A(D_A, D_B)$$
$$Y_B(D_A, D_B, Z_A, Z_B) = Y_B(D_A, D_B)$$

**A2** Random assignment:

$$[Y_A(D_A, D_B), Y_B(D_A, D_B), D_A(Z_A, Z_B), D_B(Z_A, Z_B)] \perp\!\!\!\perp [Z_A, Z_B]$$

**A3** i.i.d. assignment:  $Z_A \perp\!\!\!\perp Z_B$  and  $\mathbb{P}(Z_A = 1) = \mathbb{P}(Z_B = 1) \equiv P \in (0, 1)$

**A4** *One-sided noncompliance*:

$$D_A(Z_A = 0, Z_B) = D_B(Z_A, Z_B = 0) = 0 \forall Z_A, Z_B$$

**A5** *Monotonicity*: 
$$D_A(Z_A = 1, Z_B = 0) \leq D_A(Z_A = 1, Z_B = 1)$$
$$D_B(Z_A = 0, Z_B = 1) \leq D_B(Z_A = 1, Z_B = 1)$$

**A6** *Invertibility*:

$$\mathbb{P}(D_A(10) = 1, D_A(11) = 1, D_B(01), D_B(11) = 1) > 0$$

where

$$D_A(10) \equiv D_A(Z_A = 1, Z_B = 0)$$

$$D_A(11) \equiv D_A(Z_A = 1, Z_B = 1)$$

$$D_B(01) \equiv D_A(Z_A = 0, Z_B = 1)$$

$$D_B(11) \equiv D_A(Z_A = 1, Z_B = 1)$$

### 3.4.1 Exclusion

The outcome variable is not influenced directly by the assignment, only by the actual treatment participation. That is, the fact itself that the unit is instructed to take the treatment has no effect whatsoever on the outcome of interest – it is only the fact whether the pair members participate in the treatments which may influence the outcome.

### 3.4.2 Random Assignment

The assignment must be independent of the outcomes and the actual treatment participation. This rules out that the experimenter encourages those units to take the treatment who is more likely to (i) benefit from it (e.g. will have higher/lower outcome) or (ii) follow their treatment instruction (e.g. will have higher participation). Hence, the experimenter must randomise the treatment assignment.

### 3.4.3 i.i.d. Assignment

Every single *unit* in the sample is independently encouraged, with the same probability, to take his/her appropriate member-specific treatment. As a result, the experimenter does not have to take into account who is the pair of whom during the assignment procedure; the only thing has to be tracked is whether the unit is member-A or member-B in the pair in which he/she is.

### 3.4.4 One-sided Noncompliance

Whenever a unit is not encouraged to take the treatment, he/she is enforcably excluded from participation, i.e. has no access to the treatment.

### 3.4.5 Monotonicity

If a member is willing to take the treatment when his/her partner is excluded from doing so, he/she must also take the treatment whenever they are both encouraged to take the treatment.

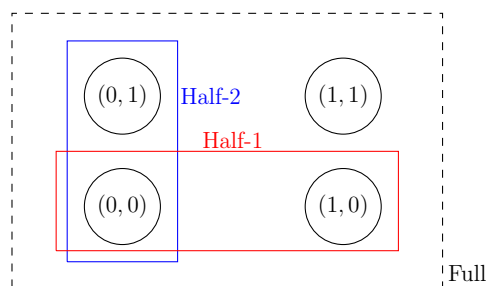
### 3.4.6 Invertibility

There must be such member-A and member-B units in the population who participate when he/she is encouraged but his/her pair is excluded from the treatment. In fact, these units form the complier subpopulations for whom the average effects are identified (see later). Intuitively, if there are no units who respect their own encouragement, ignoring their partner's, we cannot identify and estimate the effects.

### 3.5. Half-sample Identification Strategy

Having discussed the assumptions, we can turn to the Half-sample Identification Strategy, which is costly as we throw away half of the sample. Figure 1 illustrates how the half-sample and the full-sample strategies relate to one another.

**Figure 1: Half- and Full-sample strategies, with  $(Z_A, Z_B)$  values in the nodes**



As discussed earlier, the sample can be split up into groups  $\mathcal{G}_{00}, \mathcal{G}_{10}, \mathcal{G}_{01}, \mathcal{G}_{11}$  based on the sample values of  $(Z_A, Z_B)$ . In the Full-sample Identification Strategy we use the observations in all four groups. In the Half-sample Identification Strategy we use observations in either groups  $\mathcal{G}_{00}$  and  $\mathcal{G}_{10}$  (Half-1), or  $\mathcal{G}_{00}$  and  $\mathcal{G}_{01}$  (Half-2). In our example, these are the pairs: Half-1: (no pill, no education) and (pill, no education), or Half-2: (no pill, no education) and (no pill, education). That is, when one member is excluded from the treatment, so there are exactly two half-sample strategies.

What the half-sample strategy exploits is *A4 One-sided noncompliance*. Because under Identifying Assumptions *A1 Exclusion* and *A4 One-sided noncompliance*, the observable outcome variable can be written as

Half-1

$$\begin{aligned}
 Y_A^{\text{Half-1}} &= D_A Y_A(10) + (1 - D_A) Y_A(00) \\
 &= Y_A(00) + (Y_A(10) - Y_A(00)) D_A \\
 Y_B^{\text{Half-1}} &= D_A Y_B(10) + (1 - D_A) Y_B(00) \\
 &= Y_B(00) + (Y_B(10) - Y_B(00)) D_A
 \end{aligned}$$

Half-2

$$\begin{aligned}
Y_A^{\text{Half-2}} &= D_B Y_A(01) + (1 - D_B) Y_A(00) \\
&= Y_A(00) + (Y_A(01) - Y_A(00)) D_B \\
Y_B^{\text{Half-2}} &= D_B Y_B(01) + (1 - D_B) Y_B(00) \\
&= Y_B(00) + (Y_B(01) - Y_B(00)) D_B.
\end{aligned}$$

From this follows Theorem 1.

**Theorem 1** (Half-sample Identification). *Applying the IV method of Imbens and Angrist (1994) in each of the cases (Half-1 and Half-2) separately (that is, instrumenting  $D_A$  with  $Z_A$  and  $D_B$  with  $Z_B$ ), the Paired  $2 \times 2$  Factorial Design and Identifying Assumptions A4 – A6 are sufficient to identify:*

- [1] Average baseline outcomes:  $\mathbb{E}[Y_A(00)]$  and  $\mathbb{E}[Y_B(00)]$
- [2] Average effect of own treatment for compliers (ATEO):
  - $\mathbb{E}[Y_A(10) - Y_A(00) \mid D_A(10) = 1, D_A(11) = 1]$
  - $\mathbb{E}[Y_B(01) - Y_B(00) \mid D_B(01) = 1, D_B(11) = 1]$
- [3] Average effect of partner's treatment for those with complier partner (ATEP):
  - $\mathbb{E}[Y_A(01) - Y_A(00) \mid D_B(01) = 1, D_B(11) = 1]$
  - $\mathbb{E}[Y_B(10) - Y_B(00) \mid D_A(10) = 1, D_A(11) = 1].$

*Proof: see Appendix B.*

According to Theorem 1 we can identify the following averages for the member-A:

- [1]  $\mathbb{E}[Y_A(00)]$  which is the average baseline level for member-A in the whole population, the expected outcome when none of the members is treated. In our example (member-A is the depressed), this is the average outcome of the depressed member when he/she does not take the antidepressant, and nor does his/her spouse participate in the educational program.
- [2]  $\mathbb{E}[Y_A(10) - Y_A(00) \mid D_A(11) = 1, D_A(10) = 1]$  which is the average treatment effect of his/her own (member-specific) treatment on member-A, in the complier subpopulation of member-A's, i.e. those who respect their own treatment encouragement regardless of their partner's access to treatment. In our example, this is the average treatment effect of the antidepressant (given to the depressed member) on the depressed member when the non-depressed has no access to the educational program for the (sub)population of those depressed people who are willing to take the antidepressant regardless of their spouse's access to the educational program.
- [3]  $\mathbb{E}[Y_A(01) - Y_A(00) \mid D_B(11) = 1, D_B(01) = 1]$  which is the average treatment effect of his/her partner's (member-specific) treatment on member-A, in the subpopulation



of member-A's who have complier partners, i.e. those whose partner respect their treatment encouragement, regardless of member-A's access to treatment. In our example, this is the average treatment effect of the educational program (given to the non-depressed member) on the depressed member when the depressed has no access to the antidepressant for the (sub)population of those depressed people who have partners that are willing to take the educational program regardless of the depressed's access to the antidepressant.

Similarly, for member-B:

- [1]  $\mathbb{E}[Y_B(00)]$  which is the average baseline level for member-B in the whole population, the expected outcome when none of the members is treated. In our example, this is the average outcome of the non-depressed member when neither him/her nor his/her spouse receives their corresponding treatment.
- [2]  $\mathbb{E}[Y_B(01) - Y_B(00) \mid D_B(11) = 1, D_B(01) = 1]$  which is the average treatment effect of his/her own (member-specific) treatment on member-B, in the complier subpopulation of member-B's, i.e. those who respect their own treatment encouragement regardless of their partner's access to treatment. In our example, this is the average treatment effect of the educational program (given to the non-depressed member) on the non-depressed member when the depressed has no access to the antidepressant for the (sub)population of those non-depressed people who are willing to take the educational program regardless of their spouse's access to the antidepressant.
- [3]  $\mathbb{E}[Y_B(10) - Y_B(00) \mid D_A(11) = 1, D_A(10) = 1]$  which is the average treatment effect of his/her partner's (member-specific) treatment on member-B, in the subpopulation of member-B's who have complier partners, i.e. those whose partner respect their treatment encouragement, regardless of member-B's access to treatment. In our example, this is the average treatment effect of the antidepressant (given to the depressed member) on the non-depressed member when the non-depressed has no access to the educational program for the (sub)population of those non-depressed people who have partners that are willing to take the antidepressant regardless of the non-depressed's access to the educational program.

There are only two half-sample strategies in Figure 1. Why only two if there are more possibilities? In the same spirit, it is tempting to pick any other two groups and compare the outcomes similarly. The problem is that comparison will not be simpler than in the full-sample case and/or there is not enough variation in the instrument(s) (treatment encouragement(s)) or the correlation between the encouragement and the participation is insufficient (singular matrix). I do not provide details on this here – it can be shown by going through the initial steps of the proof of Theorem 2 in the Appendix.

### 3.6. Full-sample Identification Strategy

The Full-sample Identification Strategy, the main theoretical advancement of my work, is also illustrated in Figure 1. The name is self-explanatory: the strategy has the advantage over the half-sample one that we use information present in the whole sample without requiring additional assumptions.

The experimenter can refer to Theorem 2 for the identifiable treatment effects. The identified effects [1] – [3] are identical to those in Theorem 1 and are described in Section 3.5. What, in addition, is identified is effect [4]. This is a spillover-like effect which is quite hard to interpret (in words).

**Theorem 2** (Full-sample Identification). *Let*

$$\hat{\theta}_A \equiv (n^{-1} \mathbf{Z}'_A \mathbf{D}_A)^{-1} n^{-1} \mathbf{Z}'_A \mathbf{y}_A = \left( n^{-1} \sum_{i=1}^n z_{Ai} \mathbf{d}'_{Ai} \right)^{-1} n^{-1} \sum_{i=1}^n z_{Ai} y_{Ai}$$

$$\hat{\theta}_B \equiv (n^{-1} \mathbf{Z}'_B \mathbf{D}_B)^{-1} n^{-1} \mathbf{Z}'_B \mathbf{y}_B = \left( n^{-1} \sum_{i=1}^n z_{Bi} \mathbf{d}'_{Bi} \right)^{-1} n^{-1} \sum_{i=1}^n z_{Bi} y_{Bi}.$$

*Then the Paired 2 × 2 Factorial Design and Identifying Assumptions A4 – A are sufficient to identify:*

[1] *Average baseline outcome:*  $\mathbb{E}[Y_A(00)]$

[2] *Average effect of own treatment for compliers (ATEO):*

$$\mathbb{E}[Y_A(10) - Y_A(00) \mid D_A(11) = 1, D_A(10) = 1]$$

[3] *Average effect of partner's treatment for those with complier partner (ATEP):*

$$\mathbb{E}[Y_A(01) - Y_A(00) \mid D_B(11) = 1, D_B(01) = 1]$$

[4] *Spillover-like effect:*

$$\begin{aligned} & \mathbb{E}[Y_A(11) - Y_A(10) - [Y_A(10) - Y_A(00)] \mid D_A(11) = 1, D_B(11) = 1] \\ & + \frac{\tilde{r}}{\tilde{q}} \{ \mathbb{E}[Y_A(10) - Y_A(00) \mid D_A(11) = 1, D_A(10) = 0] \\ & - \mathbb{E}[Y_A(10) - Y_A(00) \mid D_A(11) = 1, D_A(10) = 1] \} \\ & + \frac{\tilde{\mu}\tilde{r}}{\tilde{q}} \{ \mathbb{E}[Y_A(01) - Y_A(00) \mid D_B(11) = 1, D_B(01) = 0] \\ & - \mathbb{E}[Y_A(01) - Y_A(00) \mid D_B(11) = 1, D_B(01) = 1] \} \end{aligned}$$

*as the probability limit  $\text{plim } \hat{\theta}_A$ , and:*

[1] *Average baseline outcome:*  $\mathbb{E}[Y_B(00)]$

[2] *Average effect of own treatment for compliers (ATEO):*

$$\mathbb{E}[Y_B(01) - Y_B(00) \mid D_B(11) = 1, D_B(01) = 1]$$

[3] *Average effect of partner's treatment for those with complier partner (ATEP):*

$$\mathbb{E}[Y_B(10) - Y_B(00) \mid D_A(11) = 1, D_A(10) = 1]$$

[4] *Spillover-like effect:*

$$\begin{aligned} & \mathbb{E}[Y_B(11) - Y_B(01) - [Y_B(01) - Y_B(00)] \mid D_B(11) = 1, D_A(11) = 1] \\ & + \frac{\tilde{\mu}\tilde{r}}{\bar{q}} \{ \mathbb{E}[Y_B(01) - Y_B(00) \mid D_B(11) = 1, D_B(01) = 0] \\ & - \mathbb{E}[Y_B(01) - Y_B(00) \mid D_B(11) = 1, D_B(01) = 1] \} \\ & + \frac{\tilde{r}}{\bar{q}} \{ \mathbb{E}[Y_B(10) - Y_B(00) \mid D_A(11) = 1, D_A(10) = 0] \\ & - \mathbb{E}[Y_B(10) - Y_B(00) \mid D_A(11) = 1, D_A(10) = 1] \} \end{aligned}$$

as the probability limit  $\text{plim } \hat{\theta}_B$ , where

$$\begin{aligned} \tilde{r} & \equiv \mathbb{P}(D_A(11) = 1, D_A(10) = 0) \\ \tilde{\mu} & \equiv \frac{\mathbb{P}(D_B(11) = 1, D_B(01) = 0)}{\mathbb{P}(D_A(11) = 1, D_A(10) = 0)} \\ \bar{q} & \equiv \mathbb{P}(D_A(11) = 1, D_B(11) = 1) \end{aligned}$$

*Proof: see Appendix C.*

## 4. Conclusion, Limitations and Further Research

I proposed an experimental design and two treatment effect identification strategies (based on instrumental variables) when there is interference within predefined pairs, made up by distinguishable members (e.g. healthy-ill members), and there is imperfect compliance, probably of endogenous nature, with treatment encouragement. As opposed to previous works addressing these two issues, my strategies allow for within-pair discussion and coordination of treatment participation. Both strategies are suitable for identifying the average treatment effects of two, potentially member-specific, binary treatments on two, potentially member-specific, outcomes belonging to the different pair members (or on one single outcome belonging to a pair) in complier subpopulations. Hence, we can identify an own effect and a partner effect for both members separately. Unfortunately, the joint effect of the two treatments is not identified.

One of the strategies uses half of the sample, while the other one uses the whole sample. As both of them requires the same assumptions to hold, the latter strategy is superior to the half-sample one and can be regarded as the main contribution of this paper.

Some drawbacks of my method are as follows. First, generalisation to more than two units in a group is hardly feasible due to the exponentially growing number of combinations. Second, working with distinguishable pair members has its advantages on one

hand (member-specific treatments and outcomes), but is restrictive on the other because pairs with indistinguishable members require a different identification strategy. To the best of my knowledge, there is no such strategy which permits coordination of treatment participation, so this remains a subject of further research.

Last, future research topics could include inference, that is, to establish the asymptotic distribution of the estimators. Besides, one could consider blocking on pretreatment variables, i.e. to implement the Paired  $2 \times 2$  Factorial Design within a block design, to reduce asymptotic variance.

## References

- Athey, S. and Imbens, G. (2016). The Econometrics of Randomized Experiments. *ArXiv e-prints*.
- Barkley, B. G., Hudgens, M. G., Clemens, J. D., Ali, M., and Emch, M. E. (2017). Causal Inference from Observational Studies with Clustered Interference. *ArXiv e-prints*.
- Blackwell, M. (2017). Instrumental Variable Methods for Conditional Effects and Causal Interaction in Voter Mobilization Experiments. *Journal of the American Statistical Association*, 112(518):590–599.
- Cheng, C.-S. (2013). *Theory of Factorial Design: Single- and Multi-Stratum Experiments*. Chapman & Hall/CRC Monographs on Statistics & Applied Probability. Taylor and Francis, Hoboken, NJ.
- Dai, B., Ding, S., and Wahba, G. (2012). Multivariate Bernoulli distribution. *ArXiv e-prints*.
- Halloran, M. E. and Struchiner, C. J. (1991). Study Designs for Dependent Happenings. *Epidemiology*, 2(5):331–338.
- Halloran, M. E. and Struchiner, C. J. (1995). Causal Inference in Infectious Diseases. *Epidemiology*, 6(2):142–151.
- Hudgens, M. G. and Halloran, M. E. (2008). Toward Causal Inference with Interference. *Journal of the American Statistical Association*, 103(482):832–842.
- Imbens, G. W. and Angrist, J. D. (1994). Identification and Estimation of Local Average Treatment Effects. *Econometrica*, 62(2):467–475.
- Kang, H. and Imbens, G. (2016). Peer Encouragement Designs in Causal Inference with Partial Interference and Identification of Local Average Network Effects. *ArXiv e-prints*.
- Merrill, P. D. and McClure, L. A. (2015). Dichotomizing Partial Compliance and Increased Participant Burden in Factorial Designs: The Performance of Four Noncompliance Methods. 16:523.
- Rosenbaum, P. R. (2007). Interference Between Units in Randomized Experiments. *Journal of the American Statistical Association*, 102(477):191–200.
- Rubin, D. B. (1974). Estimating Causal Effects of Treatments in Randomized and Nonrandomized Studies. *Journal of Educational Psychology*, 66(5):688.

- Rubin, D. B. (1980). Randomization Analysis of Experimental Data: The Fisher Randomization Test Comment. *Journal of the American Statistical Association*, 75(371):591–593.
- Saul, B. and Hudgens, M. (2017). A Recipe for Inference: Start with Causal Inference. Add Interference. Mix Well with R. *Journal of Statistical Software, Articles*, 82(2):1–21.
- Sobel, M. E. (2006). What Do Randomized Studies of Housing Mobility Demonstrate? Causal Inference in the Face of Interference. *Journal of the American Statistical Association*, 101(476):1398–1407.
- Tchetgen, E. J. T. and VanderWeele, T. J. (2012). On Causal Inference in the Presence of Interference. *Statistical Methods in Medical Research*, 21(1):55–75. PMID: 21068053.
- VanderWeele, T. J., Tchetgen, E. J. T., and Halloran, M. E. (2014). Interference and Sensitivity Analysis. *Statistical Science: A Review Journal of the Institute of Mathematical Statistics*, 29(4):687.

# Appendices

## A. Distribution of Treatment Participation

The purpose of examining the distribution of the treatment participation is two-fold. First, to see what the Identifying Assumptions  $A5 - A6$  require from the distribution to exhibit (Appendix A.1). Second, to derive implications of Identifying Assumptions which are later used in the proof of Theorem 1 and 2 (Appendix A.2 and ??).

As previously, I use the abbreviations  $D_A(10)$  for  $D_A(Z_A = 1, Z_B = 0)$ ,  $D_A(01)$  for  $D_A(Z_A = 0, Z_B = 1)$ ,  $D_B(10)$  for  $D_B(Z_A = 1, Z_B = 0)$ ,  $D_B(01)$  for  $D_B(Z_A = 0, Z_B = 1)$  and so on. During the analysis,  $A4$  *One sided-noncompliance is assumed to hold*, leading to the degenerate random variables  $D_A(00) = D_A(01) = D_B(00) = D_B(10) = 0$ . Hence the distribution boils down to that of  $\mathcal{D} \equiv [D_A(10), D_A(11), D_B(01), D_B(11)]'$ . This is a 4-variate Bernoulli distribution supported on  $\{0, 1\}^4$  and can be characterised by 16 non-negative parameters (Dai et al., 2012), each of which indicates the joint probability of a certain combination of zeros and ones:

$$\begin{aligned}
 p_{0000} &\equiv \mathbb{P}(D_A(10) = 0, D_A(11) = 0, D_B(01) = 0, D_B(11) = 0) \\
 p_{1000} &\equiv \mathbb{P}(D_A(10) = 1, D_A(11) = 0, D_B(01) = 0, D_B(11) = 0) \\
 p_{0100} &\equiv \mathbb{P}(D_A(10) = 0, D_A(11) = 1, D_B(01) = 0, D_B(11) = 0) \\
 &\vdots \\
 p_{0001} &\equiv \mathbb{P}(D_A(10) = 0, D_A(11) = 0, D_B(01) = 0, D_B(11) = 1) \\
 p_{1111} &\equiv \mathbb{P}(D_A(10) = 1, D_A(11) = 1, D_B(01) = 1, D_B(11) = 1).
 \end{aligned}$$

Let  $p$ 's be denoted by  $\mathcal{P}$ , then the probability mass function of  $\mathcal{D}$  is

$$\begin{aligned}
 f_{\mathcal{P}}(\mathcal{D}) &= p_{1111}^{D_A(10)D_A(11)D_B(01)D_B(11)} p_{0111}^{[1-D_A(01)][D_A(11)D_B(01)D_B(11)]} \times \dots \\
 &\quad \times p_{0000}^{[1-D_A(10)][1-D_A(11)][1-D_B(01)][1-D_B(11)]}.
 \end{aligned}$$

### A.1. Conditions for Identifying Assumptions

*A5 Monotonicity.* Monotonicity requires  $D_A(10) \leq D_A(11)$  and  $D_B(01) \leq D_B(11)$ . Hence whenever  $D_A(10) > D_A(11)$  or  $D_B(01) > D_B(11)$  we need  $f_{\mathcal{P}}(\mathcal{D}) = 0$  which is met by  $p_{10..} = p_{..10} = 0$ , or for more explicit form see Lemma 1.

**Lemma 1** (Conditions for  $A5$  Monotonicity). *For Identifying Assumption  $A5$  Monotonicity to hold the condition*

$$p_{1000} = p_{0010} = p_{0110} = p_{1010} = p_{1001} = p_{1011} = p_{1110} = 0.$$

is necessary and sufficient.

*A6 Invertibility.* Having clarified the implications of *A5*, let us turn our attention to *A6 Invertibility* which has

$$\mathbb{P}(D_A(1,0) = 1, D_A(1,1) = 1, D_B(0,1), D_B(1,1) = 1) \neq 0$$

from which follows Lemma 2.

**Lemma 2** (Conditions for *A6 Invertibility*). *For Identifying Assumptions A6 Invertibility to hold the condition*

$$p_{1111} > 0$$

is necessary and sufficient.

*Proof of Lemma 2.* Follows directly from the definition of  $p_{1111}$  ■

## A.2. Implications of Identifying Assumptions

*A5 Monotonicity* has the implications presented in Lemma 3, which play a role in the proof of Theorem 2.

**Lemma 3** (Implications of *A5 Monotonicity*). *A5 Monotonicity implies that*

$$\mathbb{P}(D_A(10) = 1) = \mathbb{P}(D_A(10) = 1, D_A(11) = 1)$$

$$\mathbb{P}(D_B(01) = 1) = \mathbb{P}(D_B(01) = 1, D_B(11) = 1)$$

which in turn implies that

$$\mathbb{P}(D_A(10) = 1) = \mathbb{P}(D_A(11) = 1) - \mathbb{P}(D_A(10) = 0, D_A(11) = 1)$$

$$\mathbb{P}(D_B(01) = 1) = \mathbb{P}(D_B(11) = 1) - \mathbb{P}(D_B(01) = 0, D_B(11) = 1)$$

*Proof of Lemma 3.* By the law of total probability we have

$$\mathbb{P}(D_A(10) = 1) = \mathbb{P}(D_A(10) = 1, D_A(11) = 0) + \mathbb{P}(D_A(10) = 1, D_A(11) = 1)$$

$$\mathbb{P}(D_B(01) = 1) = \mathbb{P}(D_B(01) = 1, D_B(11) = 0) + \mathbb{P}(D_B(01) = 1, D_B(11) = 1)$$

where by *A5 Monotonicity*  $\mathbb{P}(D_A(10) = 1, D_A(11) = 0) = \mathbb{P}(D_B(01) = 1, D_B(11) = 0) = 0$ , so the first claim is proven. To prove the implication of this, note that by the law of

total probability

$$\mathbb{P}(D_A(11) = 1) = \mathbb{P}(D_A(10) = 0, D_A(11) = 1) + \mathbb{P}(D_A(10) = 1, D_A(11) = 1)$$

$$\mathbb{P}(D_B(11) = 1) = \mathbb{P}(D_B(01) = 0, D_B(11) = 1) + \mathbb{P}(D_B(01) = 1, D_B(11) = 1)$$

which implies that

$$\mathbb{P}(D_A(11) = 1) - \mathbb{P}(D_A(10) = 0, D_A(11) = 1) = \mathbb{P}(D_A(10) = 1, D_A(11) = 1)$$

$$\mathbb{P}(D_B(11) = 1) - \mathbb{P}(D_B(01) = 0, D_B(11) = 1) = \mathbb{P}(D_B(01) = 1, D_B(11) = 1)$$

By the first claim of this lemma,  $\mathbb{P}(D_A(10) = 1, D_A(11) = 1) = \mathbb{P}(D_A(10) = 1)$  and  $\mathbb{P}(D_B(01) = 1, D_B(11) = 1) = \mathbb{P}(D_B(01) = 1)$  so

$$\mathbb{P}(D_A(11) = 1) - \mathbb{P}(D_A(10) = 0, D_A(11) = 1) = \mathbb{P}(D_A(10) = 1)$$

$$\mathbb{P}(D_B(11) = 1) - \mathbb{P}(D_B(01) = 0, D_B(11) = 1) = \mathbb{P}(D_B(01) = 1) \blacksquare$$

A6 Invertibility has three implications, which ensures the existence of the inverse matrices in Theorems 1 and 2.

**Lemma 4** (Implications of A6 Invertibility). *A6 Invertibility, i.e. the condition  $p_{1111} > 0$ , implies that all inequalities*

$$\mathbb{P}(D_A(10) = 1, D_A(11) = 1) \neq 0$$

$$\mathbb{P}(D_B(01) = 1, D_B(11) = 1) \neq 0$$

$$\mathbb{P}(D_A(11) = 1, D_B(11) = 1) \neq 0$$

*hold.*



*Proof of Lemma 4.* By the law of total probability we have

$$\begin{aligned}
\mathbb{P}(D_A(10) = 1, D_A(11) = 1) &= \mathbb{P}(D_A(10) = 1, D_A(11) = 1, D_B(01) = 0, D_B(11) = 0) \\
&\quad + \mathbb{P}(D_A(10) = 1, D_A(11) = 1, D_B(01) = 1, D_B(11) = 0) \\
&\quad + \mathbb{P}(D_A(10) = 1, D_A(11) = 1, D_B(01) = 0, D_B(11) = 1) \\
&\quad + \mathbb{P}(D_A(10) = 1, D_A(11) = 1, D_B(01) = 1, D_B(11) = 1) \\
&= p_{1100} + p_{1110} + p_{1101} + p_{1111} \\
\mathbb{P}(D_B(01) = 1, D_B(11) = 1) &= \mathbb{P}(D_A(10) = 0, D_A(11) = 0, D_B(01) = 1, D_B(11) = 1) \\
&\quad + \mathbb{P}(D_A(10) = 1, D_A(11) = 0, D_B(01) = 1, D_B(11) = 1) \\
&\quad + \mathbb{P}(D_A(10) = 0, D_A(11) = 1, D_B(01) = 1, D_B(11) = 1) \\
&\quad + \mathbb{P}(D_A(10) = 1, D_A(11) = 1, D_B(01) = 1, D_B(11) = 1) \\
&= p_{0011} + p_{1011} + p_{0111} + p_{1111} \\
\mathbb{P}(D_A(11) = 1, D_B(11) = 1) &= \mathbb{P}(D_A(01) = 0, D_A(11) = 1, D_B(01) = 0, D_B(11) = 1) \\
&\quad + \mathbb{P}(D_A(01) = 1, D_A(11) = 1, D_B(01) = 0, D_B(11) = 1) \\
&\quad + \mathbb{P}(D_A(01) = 0, D_A(11) = 1, D_B(01) = 1, D_B(11) = 1) \\
&\quad + \mathbb{P}(D_A(01) = 1, D_A(11) = 1, D_B(01) = 1, D_B(11) = 1) \\
&= p_{0101} + p_{1101} + p_{0111} + p_{1111}
\end{aligned}$$

Therefore  $p_{1111} > 0$  implies that all three probabilities are strictly positive ■

## B. Proof of Theorem 1

The proof of Theorem 1 directly follows from the results of [Imbens and Angrist \(1994\)](#). When using  $\mathcal{G}_{10}$  and  $\mathcal{G}_{00}$ , we set  $D(1) \equiv D_A(Z_A = 1, Z_B = 0)$ ,  $D(0) \equiv D_A(Z_A = 0, Z_B = 0)$ , and  $Y(1) \equiv Y_j(D_A = 1, D_B = 0)$ ,  $Y(0) \equiv Y_j(D_A = 0, D_B = 0)$  for  $j \in \{A, B\}$ . Directly applying the results of [Imbens and Angrist \(1994\)](#) gives the own effect for member-A (in the subpopulation of complier member-A's) and the partner effect for member-B (in the subpopulation with complier member-A partner):

$$\begin{aligned}
&\mathbb{E}[Y_A(D_A = 1, D_B = 0) - Y_A(D_A = 0, D_B = 0) \mid D_A(Z_A = 1, Z_B = 0) = 1] \\
&\mathbb{E}[Y_B(D_A = 1, D_B = 0) - Y_B(D_A = 0, D_B = 0) \mid D_A(Z_A = 1, Z_B = 0) = 1].
\end{aligned}$$

Now under  $A5$  Monotonicity the condition  $D_A(Z_A = 1, Z_B = 0) = 1$  is equivalent to the condition  $D_A(Z_A = 1, Z_B = 0) = 1$  and  $D_A(Z_A = 1, Z_B = 1) = 1$ .

When using  $\mathcal{G}_{01}$  and  $\mathcal{G}_{00}$ , we set  $D(1) \equiv D_B(Z_A = 0, Z_B = 1)$ ,  $D(0) \equiv D_B(Z_A = 0, Z_B = 0)$ , and  $Y(1) \equiv Y_j(D_A = 0, D_B = 1)$ ,  $Y(0) \equiv Y_j(D_A = 0, D_B = 0)$  for  $j \in \{A, B\}$ . Then we can immediately apply the proof of [Imbens and Angrist \(1994\)](#) again to obtain the

own effect for member-B (in the subpopulation of complier member-B's) and the partner effect for member-A (in the subpopulation with complier member-B partner):

$$\begin{aligned} & \mathbb{E}[Y_A(D_A = 0, D_B = 1) - Y_A(D_A = 0, D_B = 0) \mid D_B(Z_A = 0, Z_B = 1) = 1] \\ & \mathbb{E}[Y_B(D_A = 0, D_B = 1) - Y_B(D_A = 0, D_B = 0) \mid D_B(Z_A = 0, Z_B = 1) = 1]. \end{aligned}$$

Now under *A5 Monotonicity* the condition  $D_B(Z_A = 0, Z_B = 1) = 1$  is equivalent to the condition  $D_B(Z_A = 0, Z_B = 1) = 1$  and  $D_B(Z_A = 1, Z_B = 1) = 1$ .

The results of [Imbens and Angrist \(1994\)](#) apply as *A5 Monotonicity* together with *A6 Invertibility* implies that  $\mathbb{P}(D_A(10) = 1) \neq 0$  and  $\mathbb{P}(D_B(01) = 1) \neq 0$ . Thus the proof of [Theorem 1](#) is complete ■

## C. Proof of [Theorem 2](#)

### C.1. Generic Notation

To make the proof of [Theorem 2](#) generic for both member-A and member-B, I adopt a general  $j, j'$  indexing so that  $j \in \{A, B\}$  and  $j' \in \{A, B\} \setminus \{j\}$ . That is, when  $j$  is fixed at A,  $j'$  stands for member-B, and when  $j$  represents member-B,  $j'$  denotes member-A.

**Treatment encouragement.**  $Z_{ji}$  denotes the member- $j$  specific treatment encouragement and is 1 if and only if member- $j$  in pair  $i$  is encouraged to take the member- $j$  specific treatment. Similarly,  $Z_{j'i}$  denotes the member- $j'$  specific treatment encouragement in pair  $i$  and is 1 if and only if member- $j'$  in pair  $i$  is encouraged to take the member- $j'$  specific treatment. Suppose that  $j \equiv A$ , and member-A is the depressed person. Then  $Z_{ji} = 1$  means that the depressed person in pair  $i$  is encouraged to take the treatment.

**Treatment participation.** Let  $D_{ji}(Z_{ji}, Z_{j'i})$  denote the treatment participation of member- $j$  when the treatment encouragement of member- $j$  in pair  $i$  is  $Z_{ji}$ , while that of member- $j'$  is  $Z_{j'i}$ . It is again important to pay attention to the order of the arguments:  $D_{ji}(01)$  is the treatment participation of member- $j$ , and  $D_{j'i}(01)$  is the treatment participation of member- $j'$  in pair  $i$  when member- $j$  is not encouraged to take his/her treatment, while member- $j'$  is encouraged (i.e. the first argument always stands for member- $j$ ). Then the observable treatment participation is written as

$$\begin{aligned} D_{ji} = & Z_{ji}Z_{j'i}D_{ji}(11) + Z_{ji}(1 - Z_{j'i})D_{ji}(10) \\ & + (1 - Z_{ji})Z_{j'i}D_{ji}(01) + (1 - Z_{ji})(1 - Z_{j'i})D_{ji}(00) \end{aligned} \quad (5)$$

$$\begin{aligned} D_{j'i} = & Z_{ji}Z_{j'i}D_{j'i}(11) + Z_{ji}(1 - Z_{j'i})D_{j'i}(10) \\ & + (1 - Z_{ji})Z_{j'i}D_{j'i}(01) + (1 - Z_{ji})(1 - Z_{j'i})D_{j'i}(00). \end{aligned} \quad (6)$$

Suppose that  $j \equiv A$ , and member-A is the depressed person. Then by not giving the depressed person the antidepressant ( $Z_{ij} = 0$ ) while encouraging his/her spouse to enroll

to the educational program ( $Z_{ij'} = 1$ ), we observe  $D_{ij}(01)$  and  $D_{ij'}(01)$  telling us whether the depressed member takes the pill and whether the non-depressed one enrolls to the program given the encouragement (no pill, education program).

**Potential outcomes.** There are four potential outcomes in pair  $i$  for member- $j$ :  $Y_{ji}(D_{ji}, D_{j'i}) \in \mathbb{R}^1$  and four for member- $j'$ :  $Y_{j'i}(D_{ji}, D_{j'i}) \in \mathbb{R}^1$ . Similarly to the treatment encouragement, the observable one of the four is written as

$$\begin{aligned} Y_{ji} = & D_{ji}D_{j'i}Y_{ji}(11) + D_{ji}(1 - D_{j'i})Y_{ji}(10) \\ & + (1 - D_{ji})D_{j'i}Y_{ji}(01) + (1 - D_{ji})(1 - D_{j'i})Y_{ji}(00) \end{aligned} \quad (7)$$

$$\begin{aligned} Y_{j'i} = & D_{ji}D_{j'i}Y_{j'i}(11) + D_{ji}(1 - D_{j'i})Y_{j'i}(10) \\ & + (1 - D_{ji})D_{j'i}Y_{j'i}(01) + (1 - D_{ji})(1 - D_{j'i})Y_{j'i}(00), \end{aligned} \quad (8)$$

where the binary actual treatment participation  $(D_{ji}, D_{j'i})$  selects the observable one. In our example ( $j \equiv A =$  depressed member),  $Y_{ji}(10)$  and  $Y_{j'i}(10)$  are the outcomes of the pair members when the depressed actually takes the antidepressant and his/her spouse is not enrolled to the educational program;  $Y_{ji}(11)$  and  $Y_{j'i}(11)$  are the outcomes when the depressed takes the pill and the non-depressed is enrolled to the program and so on.

**Variables in the sample.** Let  $y_{ji}$  denote the sample analogue of the observable outcome for member- $j$ ; let  $\mathbf{d}_{ji} \equiv [1, d_{ji}, d_{j'i}, d_{ji}d_{j'i}]' \in \{0, 1\}^{4 \times 1}$ , where  $d_{ji}$  is the sample-realisation of treatment participation of member- $j$  in pair  $i$ ; and last let  $\mathbf{z}_{ji} \equiv [1, z_{ji}, z_{j'i}, z_{ji}z_{j'i}]' \in \{0, 1\}^{4 \times 1}$ , where  $z_{ji}$  is the sample-realisation of the treatment encouragement of member- $j$  in pair  $i$ . Broadcasting them to  $\mathbf{y}_j \equiv [y_{j1}, \dots, y_{jn}]' \in \mathbb{R}^{n \times 1}$ ,  $\mathbf{D}'_j \equiv [\mathbf{d}_{j1}, \dots, \mathbf{d}_{jn}] \in \{0, 1\}^{4 \times n}$ ,  $\mathbf{Z}'_j \equiv [\mathbf{z}_{j1}, \dots, \mathbf{z}_{jn}] \in \{0, 1\}^{4 \times n}$  facilitates more compact notation.

## C.2. Statement of Theorem 2

Theorem 2 states that the certain treatment effects are identified as  $\text{plim } \hat{\boldsymbol{\theta}}_j$  if Identifying Assumptions are met, where the IV-based estimator is

$$\begin{aligned} \hat{\boldsymbol{\theta}}_j & \equiv (n^{-1} \mathbf{Z}'_j \mathbf{D}_j)^{-1} n^{-1} \mathbf{Z}'_j \mathbf{y}_j \\ & = \left( n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} \mathbf{d}'_{ji} \right)^{-1} n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} y_{ji} \\ \mathbf{d}_{ji} & \equiv [1, d_{ji}, d_{j'i}, d_{ji}d_{j'i}]' \in \{0, 1\}^{4 \times 1} \\ \mathbf{z}_{ji} & \equiv [1, z_{ji}, z_{j'i}, z_{ji}z_{j'i}]' \in \{0, 1\}^{4 \times 1} \end{aligned}$$

for  $j \in \{A, B\}$  and  $j' \in \{A, B\} \setminus \{j\}$ . To prove this, we examine  $\text{plim } \hat{\boldsymbol{\theta}}_j$ , i.e. the probability limit of  $\hat{\boldsymbol{\theta}}_j$ , or equivalently  $\boldsymbol{\theta}_j : \lim_{n \rightarrow \infty} \mathbb{P} \left( \|\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j\|_2^2 > \varepsilon \right) = 0$  for any  $\varepsilon > 0$ . By the continuous mapping property of the probability limit and by the Weak Law of Large

Numbers for i.i.d. data (as the data across pairs are i.i.d.), we have

$$\begin{aligned}
\text{plim } \hat{\theta}_j &= \text{plim} \left( n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} \mathbf{d}'_{ji} \right)^{-1} n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} y_{ji} \\
&= \left( \text{plim } n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} \mathbf{d}'_{ji} \right)^{-1} \text{plim } n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} y_{ji} \\
&= \mathbb{E} [\mathbf{z}_{ji} \mathbf{d}'_{ji}]^{-1} \mathbb{E} [\mathbf{z}_{ji} y_{ji}] \\
&= \mathbb{E} \left[ \begin{bmatrix} 1 \\ Z_j \\ Z_{j'} \\ Z_j Z_{j'} \end{bmatrix} \begin{bmatrix} 1 & D_j & D_{j'} & D_j D_{j'} \end{bmatrix} \right]^{-1} \mathbb{E} \left[ \begin{bmatrix} 1 \\ Z_j \\ Z_{j'} \\ Z_j Z_{j'} \end{bmatrix} Y_j \right] \\
&\equiv \mathbf{M}^{-1} \mathbf{v}. \tag{9}
\end{aligned}$$

Thus we need to compute the appropriate expectations and calculate the product. The proof is organised into sections, so that in Appendix C.3  $\mathbf{M}$ , and in Appendix C.4  $\mathbf{v}$  are computed, and then the product  $\mathbf{M}^{-1} \mathbf{v}$  is evaluated and the results are summarised in Appendix C.6.

### C.3. Computing $\mathbf{M}$

In this section, I compute  $\mathbf{M}$ , which is defined in (9). For a start, because of A4 One-sided noncompliance ( $D_j(0, Z_{j'}) = D_{j'}(Z_j, 0) = 0 \forall Z_j, Z_{j'}$ ) we can rewrite the observable treatment participations in (10) and (12) as

$$\begin{aligned}
D_j &= Z_j Z_{j'} D_j(11) + Z_j(1 - Z_{j'}) D_j(10) \\
&\quad + (1 - Z_j) Z_{j'} D_j(01) + (1 - Z_j)(1 - Z_{j'}) D_j(00) \\
&= Z_j Z_{j'} [D_j(11) - D_j(10)] + Z_j D_j(10) \\
D_{j'} &= Z_j Z_{j'} D_{j'}(11) + Z_j(1 - Z_{j'}) D_{j'}(10) \\
&\quad + (1 - Z_j) Z_{j'} D_{j'}(01) + (1 - Z_j)(1 - Z_{j'}) D_{j'}(00) \\
&= Z_j Z_{j'} [D_{j'}(11) - D_{j'}(01)] + Z_{j'} D_{j'}(01).
\end{aligned}$$

Now we can focus on  $\mathbf{M}$ :

$$\begin{aligned}
\mathbf{M} &= \mathbb{E} \left[ \begin{bmatrix} 1 \\ Z_j \\ Z_{j'} \\ Z_j Z_{j'} \end{bmatrix} \begin{bmatrix} 1 & D_j & D_{j'} & D_j D_{j'} \end{bmatrix} \right] \\
&= \mathbb{E} \left[ \begin{bmatrix} 1 & D_j & D_{j'} & D_j D_{j'} \\ Z_j & Z_j D_j & Z_j D_{j'} & Z_j D_j D_{j'} \\ Z_{j'} & Z_{j'} D_j & Z_{j'} D_{j'} & Z_{j'} D_j D_{j'} \\ Z_j Z_{j'} & Z_j Z_{j'} D_j & Z_j Z_{j'} D_{j'} & Z_j Z_{j'} D_j D_{j'} \end{bmatrix} \right] \\
&= \begin{bmatrix} 1 & \mathbb{E}[D_j] & \mathbb{E}[D_{j'}] & \mathbb{E}[D_j D_{j'}] \\ \mathbb{E}[Z_j] & \mathbb{E}[Z_j D_j] & \mathbb{E}[Z_j D_{j'}] & \mathbb{E}[Z_j D_j D_{j'}] \\ \mathbb{E}[Z_{j'}] & \mathbb{E}[Z_{j'} D_j] & \mathbb{E}[Z_{j'} D_{j'}] & \mathbb{E}[Z_{j'} D_j D_{j'}] \\ \mathbb{E}[Z_j Z_{j'}] & \mathbb{E}[Z_j Z_{j'} D_j] & \mathbb{E}[Z_j Z_{j'} D_{j'}] & \mathbb{E}[Z_j Z_{j'} D_j D_{j'}] \end{bmatrix}.
\end{aligned}$$

In the following, these expectations are computed row-by-row, one-by-one. At each element I use: A2 Random assignment to factor the expectation of the product of the assignment(s) and participation(s) into the product of the expected assignment(s) and expected participation(s); A3 i.i.d. assignment to have  $\mathbb{E}[Z_j Z_{j'} X] = \mathbb{E}[Z_j] \mathbb{E}[Z_{j'}] \mathbb{E}[X] = PP \mathbb{E}[X] = P^2 \mathbb{E}[X]$  for any random variable  $X$  which is not  $Z_j$  or  $Z_{j'}$ . Also note that for any binary random variable  $W$  we have  $\mathbb{E}[W] = \mathbb{P}(W = 1)$  and that  $W^2 = W$ . Let us proceed with the expectations.

**Row of  $\mathbf{M}$ : 1st**

$$\begin{aligned}
\mathbb{E}[D_j] &= \mathbb{E}[Z_j Z_{j'} [D_j(11) - D_j(10)] + Z_j D_j(10)] \\
&= P^2 \mathbb{E}[D_j(11) - D_j(10)] + P \mathbb{E}[D_j(10)] \\
&= P^2 \mathbb{E}[D_j(11) - D_j(10)] + P \mathbb{P}(D_j(10) = 1) \quad \text{A5} \implies \\
&= P^2 [0 \mathbb{P}(D_j(11) = 0, D_j(10) = 0) + 1 \mathbb{P}(D_j(11) = 1, D_j(10) = 0) \\
&\quad + 0 \mathbb{P}(D_j(11) = 1, D_j(10) = 1)] + P \mathbb{P}(D_j(10) = 1) \\
&= P^2 \mathbb{P}(D_j(11) = 1, D_j(10) = 0) + P \mathbb{P}(D_j(10) = 1). \tag{10}
\end{aligned}$$

Similarly for  $D_{j'}$

$$\mathbb{E}[D_{j'}] = \mathbb{E}[Z_j Z_{j'} [D_{j'}(11) - D_{j'}(01)] + Z_{j'} D_{j'}(01)] \tag{11}$$

$$= P^2 \mathbb{P}(D_{j'}(11) = 1, D_{j'}(01) = 0) + P \mathbb{P}(D_{j'}(01) = 1). \tag{12}$$

Then

$$\begin{aligned}
D_j D_{j'} &= \{Z_j Z_{j'} [D_j(11) - D_j(10)] + Z_j D_j(10)\} \\
&\quad \times \{Z_j Z_{j'} [D_{j'}(11) - D_{j'}(01)] + Z_{j'} D_{j'}(01)\} \\
&= Z_j^2 Z_{j'}^2 [D_j(11) - D_j(10)] [D_{j'}(11) - D_{j'}(01)] + \\
&\quad + Z_j Z_{j'}^2 [D_j(11) - D_j(10)] D_{j'}(01) \\
&\quad + Z_j^2 Z_{j'} [D_{j'}(11) - D_{j'}(01)] D_j(10) \\
&\quad + Z_j Z_{j'} D_j(10) D_{j'}(01)
\end{aligned} \tag{13}$$

$$Z \in \{0, 1\} \implies$$

$$\begin{aligned}
\mathbb{E} [D_j D_{j'}] &= P^2 \{ \mathbb{E} [(D_j(11) - D_j(10))(D_{j'}(11) - D_{j'}(01))] \\
&\quad + \mathbb{E} [(D_j(11) - D_j(10))D_{j'}(01)] \\
&\quad + \mathbb{E} [(D_{j'}(11) - D_{j'}(01))D_j(10)] \\
&\quad + \mathbb{E} [D_j(10)D_{j'}(01)] \} \\
&= P^2 \{ \mathbb{E} [(D_j(11) - D_j(10))(D_{j'}(11) - D_{j'}(01))] \\
&\quad + \mathbb{E} [D_j(11)D_{j'}(01) + D_{j'}(11)D_j(10) - D_{j'}(01)D_j(10)] \} \\
&= P^2 \mathbb{E} [D_j(11)D_{j'}(11)] \\
&= P^2 \mathbb{P} (D_j(11) = 1, D_{j'}(11) = 1) .
\end{aligned}$$

**Row of  $M$ : 2nd**

$$\begin{aligned}
\mathbb{E} [Z_j D_j] &= \mathbb{E} [Z_j \{Z_j Z_{j'} [D_j(11) - D_j(10)] + Z_j D_j(10)\}] \\
&= \mathbb{E} [Z_j^2 Z_{j'} [D_j(11) - D_j(10)] + Z_j^2 D_j(10)] \quad Z_j \in \{0, 1\} \implies \\
&= \mathbb{E} [Z_j Z_{j'} [D_j(11) - D_j(10)] + Z_j D_j(10)] \\
&= \mathbb{E} [D_j] \\
&= P^2 \mathbb{P} (D_j(11) = 1, D_j(10) = 0) + P \mathbb{P} (D_j(10) = 1) .
\end{aligned}$$

$$\begin{aligned}
\mathbb{E} [Z_j D_{j'}] &= \mathbb{E} [Z_j \{Z_j Z_{j'} [D_{j'}(11) - D_{j'}(01)] + Z_{j'} D_{j'}(01)\}] \quad Z_j \in \{0, 1\} \implies \\
&= \mathbb{E} [Z_j Z_{j'} [D_{j'}(11) - D_{j'}(01)] + Z_j Z_{j'} D_{j'}(01)] \\
&= \mathbb{E} [Z_j Z_{j'} D_{j'}(11)] \\
&= P^2 \mathbb{P} (D_{j'}(11) = 1) .
\end{aligned}$$

(13) and  $Z_j \in \{0, 1\} \implies$

$$\begin{aligned}\mathbb{E}[Z_j D_j D_{j'}] &= \mathbb{E}[D_j D_{j'}] \\ &= P^2 \mathbb{P}(D_j(11) = 1, D_{j'}(11) = 1).\end{aligned}$$

**Row of  $M$ : 3rd**

$$\begin{aligned}\mathbb{E}[Z_{j'} D_j] &= \mathbb{E}[Z_{j'} \{Z_j Z_{j'} [D_j(11) - D_j(10)] + Z_j D_j(10)\}] \\ &= \mathbb{E}[Z_{j'} Z_j [D_j(11) - D_j(10)] + Z_{j'} Z_j D_j(10)] \\ &= \mathbb{E}[Z_{j'} Z_j D_j(11)] \\ &= P^2 \mathbb{P}(D_j(11) = 1).\end{aligned}$$

(12) and  $Z_{j'} \in \{0, 1\} \implies$

$$\begin{aligned}\mathbb{E}[Z_{j'} D_{j'}] &= \mathbb{E}[D_{j'}] \\ &= P^2 \mathbb{P}(D_{j'}(11) = 1, D_{j'}(01) = 0) + P \mathbb{P}(D_{j'}(01) = 1).\end{aligned}$$

(13) and  $Z_{j'} \in \{0, 1\} \implies$

$$\begin{aligned}\mathbb{E}[Z_{j'} D_j D_{j'}] &= \mathbb{E}[D_j D_{j'}] \\ &= P^2 \mathbb{P}(D_j(11) = 1, D_{j'}(11) = 1).\end{aligned}$$

**Row of  $M$ : 4th**

(10) and  $Z_j \in \{0, 1\} \implies$

$$\begin{aligned}\mathbb{E}[Z_j Z_{j'} D_j] &= \mathbb{E}[Z_{j'} D_j] \\ &= P^2 \mathbb{P}(D_j(11) = 1).\end{aligned}$$

(12) and  $Z_{j'} \in \{0, 1\} \implies$

$$\begin{aligned}\mathbb{E}[Z_j Z_{j'} D_{j'}] &= \mathbb{E}[Z_j D_{j'}] \\ &= P^2 \mathbb{P}(D_{j'}(11) = 1).\end{aligned}$$

(13) and  $Z_j, Z_{j'} \in \{0, 1\} \implies$

$$\begin{aligned}\mathbb{E}[Z_j Z_{j'} D_j D_{j'}] &= \mathbb{E}[D_j D_{j'}] \\ &= P^2 \mathbb{P}(D_j(11) = 1, D_{j'}(11) = 1).\end{aligned}$$

Next, I introduce new notation for compactness:

$$\begin{aligned}
q &\equiv \mathbb{P}(D_j(11) = 1) \\
\lambda &\equiv \frac{\mathbb{P}(D_{j'}(11) = 1)}{\mathbb{P}(D_j(11) = 1)} \implies \mathbb{P}(D_{j'}(11) = 1) = \lambda q \\
r &\equiv \mathbb{P}(D_j(11) = 1, D_j(10) = 0) \\
\mu &\equiv \frac{\mathbb{P}(D_{j'}(11) = 1, D_{j'}(01) = 0)}{\mathbb{P}(D_j(11) = 1, D_j(10) = 0)} \implies \mathbb{P}(D_{j'}(11) = 1, D_{j'}(01) = 0) = \mu r \\
\bar{q} &\equiv \mathbb{P}(D_j(11) = 1, D_{j'}(11) = 1)
\end{aligned}$$

then by invoking the implications of A5 Monotonicity in Lemma 3, we can write

$$\mathbb{P}(D_j(10) = 1) = \mathbb{P}(D_j(11) = 1) - \mathbb{P}(D_j(11) = 1, D_j(10) = 0) = q - r \quad (14)$$

$$\mathbb{P}(D_{j'}(01) = 1) = \mathbb{P}(D_{j'}(11) = 1) - \mathbb{P}(D_{j'}(11) = 1, D_{j'}(01) = 0) = \lambda q - \mu r \quad (15)$$

and, because by Lemma 3  $\mathbb{P}(D_j(10) = 1) = \mathbb{P}(D_j(10) = 1, D_j(11) = 1)$  and  $\mathbb{P}(D_{j'}(01) = 1) = \mathbb{P}(D_{j'}(01) = 1, D_{j'}(11) = 1)$ ,

$$\mathbb{P}(D_j(10) = 1, D_j(11) = 1) = q - r \quad (16)$$

$$\mathbb{P}(D_{j'}(01) = 1, D_{j'}(11) = 1) = \lambda q - \mu r \quad (17)$$

Then we can compactly express the expectations of  $\mathbf{M}$  row-by-row:

**Row of  $M$ : 1st**

$$\begin{aligned}
\mathbb{E}[D_j] &= P^2 r + P(q - r) \\
\mathbb{E}[D_{j'}] &= P^2 \mu r + P(\lambda q - \mu r) \\
\mathbb{E}[D_j D_{j'}] &= P^2 \bar{q}
\end{aligned}$$

**Row of  $M$ : 2nd**

$$\begin{aligned}
\mathbb{E}[Z_j D_j] &= P^2 r + P(q - r) \\
\mathbb{E}[Z_j D_{j'}] &= P^2 \lambda q \\
\mathbb{E}[Z_j D_j D_{j'}] &= P^2 \bar{q}
\end{aligned}$$



**Row of  $M$ : 3rd**

$$\begin{aligned}\mathbb{E} [Z_{j'} D_j] &= P^2 q \\ \mathbb{E} [Z_{j'} D_{j'}] &= P^2 \mu r + P(\lambda q - \mu r) \\ \mathbb{E} [Z_{j'} D_j D_{j'}] &= P^2 \bar{q}\end{aligned}$$

**Row of  $M$ : 4th**

$$\begin{aligned}\mathbb{E} [Z_j Z_{j'} D_j] &= P^2 q \\ \mathbb{E} [Z_j Z_{j'} D_{j'}] &= P^2 \lambda q \\ \mathbb{E} [Z_j Z_{j'} D_j D_{j'}] &= P^2 \bar{q}.\end{aligned}$$

For further simplification let

$$\begin{aligned}a &\equiv P^2 r + P(q - r) \\ b &\equiv P^2 \mu r + P(\lambda q - \mu r) \\ c &\equiv P^2 \bar{q} \\ d &\equiv P^2 \lambda q \\ e &\equiv P^2 q\end{aligned}$$

so we can write

$$\mathbf{M} = \begin{bmatrix} 1 & a & b & c \\ P & a & d & c \\ P & e & b & c \\ P^2 & e & d & c \end{bmatrix}.$$

The letters  $a, b, c, d, e$  will not be used later on, their only purpose is to see how  $\mathbf{M}$  can be written in the simplest form that the reader can follow more easily `symbolic_inversion.py`.

Before moving on to  $\mathbf{v}$  in the next subsection, we have to clarify that under which conditions  $\mathbf{M}^{-1}$  exists. A square matrix is invertible if and only if its determinant is not zero, thus follow Lemma 5, 6, and 7.

**Lemma 5** (Determinant of  $\mathbf{M}$ ).

$$\det(\mathbf{M}) = P^4 \bar{q} (P - 1)^4 (q - r) (\lambda q - \mu r)$$

*Proof of Lemma 5.* See Python code `symbolic_inversion.py` ■

**Lemma 6** (Existence of  $\mathbf{M}^{-1}$ ). *Given  $P \in (0, 1)$ ,  $\mathbf{M}^{-1}$  exists if and only if all three inequalities*

$$\begin{aligned} 0 \neq q - r &= \mathbb{P}(D_j(11) = 1) - \mathbb{P}(D_j(10) = 0, D_j(11) = 1) \\ 0 \neq \lambda q - \mu r &= \mathbb{P}(D_{j'}(11) = 1) - \mathbb{P}(D_{j'}(01) = 0, D_{j'}(11) = 1) \\ 0 \neq \bar{q} &= \mathbb{P}(D_j(11) = 1, D_{j'}(11) = 1) \end{aligned}$$

*hold.*

*Proof of Lemma 6.* A square matrix is invertible if and only if its determinant is nonzero. Then the inequalities directly follow from Lemma 5 ■

**Lemma 7** (Sufficient Condition for the Existence of  $\mathbf{M}^{-1}$ ). *Given  $P \in (0, 1)$ , the condition*

$$\mathbb{P}(D_j(01) = 1, D_j(11) = 1, D_{j'}(01) = 1, D_{j'}(11) = 1) > 0$$

*implies that all three inequalities in Lemma 6 hold, therefore it is a sufficient condition for  $\mathbf{M}^{-1}$  to exist.*

*Proof of Lemma 7.* The first two inequalities in Lemma 6 can be rewritten by (16) and (17) as follows:

$$\begin{aligned} 0 \neq q - r &= \mathbb{P}(D_j(10) = 1, D_j(11) = 1) \\ 0 \neq \lambda q - \mu r &= \mathbb{P}(D_{j'}(01) = 1, D_{j'}(11) = 1) \end{aligned}$$

Now from results of Lemma 4,  $\mathbb{P}(D_j(01) = 1, D_j(11) = 1, D_{j'}(01) = 1, D_{j'}(11) = 1) > 0$  implies that the above two inequalities and the third one,  $\bar{q} = \mathbb{P}(D_j(11) = 1, D_{j'}(11) = 1) \neq 0$  hold as well ■

#### C.4. Computing $\mathbf{v}$

Having computed  $\mathbf{M}$  and checked the conditions for  $\mathbf{M}^{-1}$  to exist, now we can analyse  $\mathbf{v}$  defined in (9):

$$\begin{aligned} \mathbf{v} &= \mathbb{E} \left[ \begin{array}{c} \left[ \begin{array}{c} 1 \\ Z_j \\ Z_{j'} \\ Z_j Z_{j'} \end{array} \right] \\ Y_j \end{array} \right] \\ &= \begin{array}{c} \mathbb{E} [Y_j] \\ \mathbb{E} [Z_j Y_j] \\ \mathbb{E} [Z_{j'} Y_j] \\ \mathbb{E} [Z_j Z_{j'} Y_j] \end{array} \end{aligned}$$

This analysis is more tedious, but the procedure is the same: computing expectations one-by-one. Again, I use A2 Random assignment and A3 i.i.d. assignment throughout the computation. Before this, let us rewrite the potential outcomes in (7) as

$$\begin{aligned} Y_j &= D_j D_{j'} Y_j(11) + D_j (1 - D_{j'}) Y_j(10) \\ &\quad + (1 - D_j) D_{j'} Y_j(01) + (1 - D_j) (1 - D_{j'}) Y_j(00) \\ &= D_j D_{j'} [Y_j(11) - Y_j(01) - (Y_j(10) - Y_j(00))] \\ &\quad + D_j [Y_j(10) - Y_j(00)] + D_{j'} [Y_j(01) - Y_j(00)] + Y_j(00) \end{aligned} \tag{18}$$

and introduce  $\Delta_j(D_j, D_{j'} | \tilde{D}_j, \tilde{D}_{j'}) \equiv Y_j(D_j, D_{j'}) - Y_j(\tilde{D}_j, \tilde{D}_{j'})$ . Then we can proceed with the expectations one by one.

**Row of  $v$ : 1st,  $\mathbb{E}[Y_j]$ . (18)  $\implies$**

$$\begin{aligned}
\mathbb{E}[Y_j] &= \mathbb{E} [D_j D_{j'} \{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j \Delta_j(10|00) + D_{j'} \Delta_j(01|00) + Y_j(00)] \\
&= P^2 \mathbb{E} [D_j(11) D_{j'}(11) \{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j(11) \Delta_j(10|00) + D_{j'}(11) \Delta_j(01|00) + Y_j(00)] \\
&\quad + P(1-P) \mathbb{E} [D_j(10) D_{j'}(10) \{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j(10) \Delta_j(10|00) + D_{j'}(10) \Delta_j(01|00) + Y_j(00)] \\
&\quad + (1-P)P \mathbb{E} [D_j(01) D_{j'}(01) \{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j(01) \Delta_j(10|00) + D_{j'}(01) \Delta_j(01|00) + Y_j(00)] \\
&\quad + (1-P)^2 \mathbb{E} [D_j(00) D_{j'}(00) \{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j(00) \Delta_j(10|00) + D_{j'}(00) \Delta_j(01|00) + Y_j(00)] \tag{19}
\end{aligned}$$

(20)

Now by *A4 One-sided noncompliance* ( $D_j(0, Z_{j'}) = D_{j'}(Z_j, 0) = 0 \forall Z_j, Z_{j'}$ )

$$\begin{aligned}
\mathbb{E}[Y_j] &= P^2 \underbrace{\mathbb{E} [D_j(11) D_{j'}(11) \{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j(11) \Delta_j(10|00) + D_{j'}(11) \Delta_j(01|00) + Y_j(00)]}_{\equiv \alpha} \\
&\quad + P(1-P) \underbrace{\mathbb{E} [D_j(10) \Delta_j(10|00) + Y_j(00)]}_{\equiv \beta} \\
&\quad + (1-P)P \underbrace{\mathbb{E} [D_{j'}(01) \Delta_j(01|00) + Y_j(00)]}_{\equiv \gamma} \\
&\quad + (1-P)^2 \underbrace{\mathbb{E} [Y_j(00)]}_{\equiv \psi},
\end{aligned}$$

so

$$\begin{aligned}
\mathbb{E}[Y_j] &= P^2\alpha + P(1-P)\beta + (1-P)P\gamma + (1-P)^2\psi \\
&= P^2\alpha + P\beta - P^2\beta + P\gamma - P^2\gamma + (1-P)\psi - (P-P^2)\psi \\
&= P^2(\alpha - (\beta + \gamma)) + P(\beta + \gamma) + (1-P)\psi - (P-P^2)\psi \\
&= P^2(\alpha - (\beta + \gamma - \psi)) + P(\beta + \gamma - 2\psi) + \psi.
\end{aligned}$$

What we want to take advantage of now is *A5 Monotonicity*, which is best done by expressing as many terms as we can with  $D_j(11) - D_j(10)$  and  $D_{j'}(11) - D_{j'}(01)$ . This is why the greek letters are introduced, so we can proceed like this:

$$\begin{aligned}
\beta + \gamma - \psi &= \mathbb{E}[D_j(10)\Delta_j(10|00) + Y_j(00)] + \mathbb{E}[D_{j'}(01)\Delta_j(01|00) + Y_j(00)] - \mathbb{E}[Y_j(00)] \\
&= \mathbb{E}[D_j(10)\Delta_j(10|00) + D_{j'}(01)\Delta_j(01|00) + Y_j(00)] \tag{21} \\
\alpha - (\beta + \gamma - \psi) &= \mathbb{E}[D_j(11)D_{j'}(11)\{\Delta_j(11|01) - \Delta_j(10|00)\} + (D_j(11) - D_j(10))\Delta_j(10|00) + (D_{j'}(11) - D_{j'}(01))\Delta_j(01|00)],
\end{aligned}$$

which simplifies by *A5 Monotonicity* as follows

$$\begin{aligned}
\mathbb{E}[D_j(11)D_{j'}(11)\{\Delta_j(11|01) - \Delta_j(10|00)\}] &= \mathbb{E}[\Delta_j(11|01) - \Delta_j(10|00) \mid D_j(11) = 1, D_{j'}(11) = 1] \mathbb{P}(D_j(11) = 1, D_{j'}(11) = 1) \\
\mathbb{E}[(D_j(11) - D_j(10))\Delta_j(10|00)] &= \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \mathbb{P}(D_j(11) = 1, D_j(10) = 0) \\
\mathbb{E}[(D_{j'}(11) - D_{j'}(01))\Delta_j(01|00)] &= \mathbb{E}[\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 0] \mathbb{P}(D_{j'}(11) = 1, D_{j'}(01) = 0),
\end{aligned}$$

thus

$$\begin{aligned}\alpha - (\beta + \gamma - \psi) &= \mathbb{E} [\Delta_j(11|01) - \Delta_j(10|00) \mid D_j(11) = 1, D_{j'}(11) = 1] \mathbb{P} (D_j(11) = 1, D_{j'}(11) = 1) \\ &\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \mathbb{P} (D_j(11) = 1, D_j(10) = 0) \\ &\quad + \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(10) = 0] \mathbb{P} (D_{j'}(11) = 1, D_{j'}(01) = 0).\end{aligned}$$

Next, by (21) and the definition of  $\psi = \mathbb{E}[Y_j(00)]$ ,

$$\begin{aligned}\beta + \gamma - 2\psi &= \mathbb{E} [D_j(10)\Delta_j(10|00) + D_{j'}(01)\Delta_j(01|00)] \\ &= \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1] \mathbb{P} (D_j(10) = 1) + \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(01) = 1] \mathbb{P} (D_{j'}(01) = 1).\end{aligned}$$

Putting together the greek-letter expressions, we finally obtain the first row of  $\mathbf{v}$ :

$$\begin{aligned}\mathbb{E} [Y_j] &= P^2(\alpha - (\beta + \gamma - \psi)) + P(\beta + \gamma - 2\psi) + \psi \\ &= P^2\{\mathbb{E} [\Delta_j(11|01) - \Delta_j(10|00) \mid D_j(11) = 1, D_{j'}(11) = 1] \mathbb{P} (D_j(11) = 1, D_{j'}(11) = 1) \\ &\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \mathbb{P} (D_j(11) = 1, D_j(10) = 0) \\ &\quad + \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(10) = 0] \mathbb{P} (D_{j'}(11) = 1, D_{j'}(01) = 0)\} \\ &\quad + P\{\mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1] \mathbb{P} (D_j(10) = 1) + \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(01) = 1] \mathbb{P} (D_{j'}(01) = 1)\} \\ &\quad + \mathbb{E} [Y_j(00)].\end{aligned}$$

**Row of  $v$ : 2nd**,  $\mathbb{E}[Z_j Y_j]$ . (18) and A4 One-sided noncompliance  $\implies$

$$\begin{aligned}
\mathbb{E}[Z_j Y_j] &= \mathbb{E}[Z_j \{D_j D_{j'} \{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j \Delta_j(10|00) + D_{j'} \Delta_j(01|00) + Y_j(00)\}] \\
&= P^2 \mathbb{E}[D_j(11) D_{j'}(11) \{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j(11) \Delta_j(10|00) + D_{j'}(11) \Delta_j(01|00) + Y_j(00)] \\
&\quad + P(1-P) \mathbb{E}[D_j(10) \Delta_j(10|00) + Y_j(00)] \\
&= P^2 \mathbb{E}[D_j(11) D_{j'}(11) \{\Delta_j(11|01) - \Delta_j(10|00)\} + (D_j(11) - D_j(10)) \Delta_j(10|00) + D_{j'}(11) \Delta_j(01|00)] \\
&\quad + P \mathbb{E}[D_j(10) \Delta_j(10|00) + Y_j(00)].
\end{aligned} \tag{22}$$

Note how the last two terms in (19) are not present in (22) because of A4 One-sided noncompliance. Next, as we did for  $\mathbb{E}[Y_j]$ , using A5 Monotonicity:

$$\begin{aligned}
\mathbb{E}[Z_j Y_j] &= P^2 \{ \mathbb{E}[\Delta_j(11|01) - \Delta_j(10|00) \mid D_j(11) = 1, D_{j'}(11) = 1] \mathbb{P}(D_j(11) = 1, D_{j'}(11) = 1) \\
&\quad + \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \mathbb{P}(D_j(11) = 1, D_j(10) = 0) \\
&\quad + \mathbb{E}[\Delta_j(01|00) \mid D_{j'}(11) = 1] \mathbb{P}(D_{j'}(11) = 1) \} \\
&+ P \mathbb{E}[\Delta_j(10|00) \mid D_j(10) = 1] \mathbb{P}(D_j(10) = 1) \\
&+ P \mathbb{E}[Y_j(00)].
\end{aligned}$$

**Row of v: 3rd**,  $\mathbb{E} [Z_{j'}Y_j]$ . (18) and A4 One-sided noncompliance  $\implies$

$$\begin{aligned}
\mathbb{E} [Z_{j'}Y_j] &= \mathbb{E} [Z_{j'}\{D_j D_{j'}\{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j \Delta_j(10|00) + D_{j'} \Delta_j(01|00) + Y_j(00)\}] \\
&= P^2 \mathbb{E} [D_j(11)D_{j'}(11)\{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j(11)\Delta_j(10|00) + D_{j'}(11)\Delta_j(01|00) + Y_j(00)] \\
&\quad + (1 - P)P \mathbb{E} [D_{j'}(01)\Delta_j(01|00)] \\
&= P^2 \mathbb{E} [D_j(11)D_{j'}(11)\{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j(11)\Delta_j(10|00) + (D_{j'}(11) - D_{j'}(01))\Delta_j(01|00)] \\
&\quad + P \mathbb{E} [D_{j'}(01)\Delta_j(01|00) + Y_j(00)]
\end{aligned} \tag{23}$$

Note again how the two terms in (19) are not present in (23) because of A4 One-sided noncompliance. Next, as we did for  $\mathbb{E} [Y_j]$ , using A5 Monotonicity:

$$\begin{aligned}
\mathbb{E} [Z_{j'}Y_j] &= P^2 \{ \mathbb{E} [\Delta_j(11|01) - \Delta_j(10|00) \mid D_j(11) = 1, D_{j'}(11) = 1] \mathbb{P} (D_j(11) = 1, D_{j'}(11) = 1) \\
&\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1] \mathbb{P} (D_j(11) = 1) \\
&\quad + \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 0] \mathbb{P} (D_{j'}(11) = 1, D_{j'}(01) = 0) \} \\
&\quad + P \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(01) = 1] \mathbb{P} (D_{j'}(01) = 1) \\
&\quad + P \mathbb{E} [Y_j(00)].
\end{aligned}$$

**Row of v: 4th**,  $\mathbb{E} [Z_j Z_{j'} Y_j]$ . (18) and A4 One-sided noncompliance  $\implies$

$$\begin{aligned}
\mathbb{E} [Z_j Z_{j'} Y_j] &= \mathbb{E} [Z_j Z_{j'}\{D_j D_{j'}\{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j \Delta_j(10|00) + D_{j'} \Delta_j(01|00) + Y_j(00)\}] \\
&= P^2 \mathbb{E} [D_j(11)D_{j'}(11)\{\Delta_j(11|01) - \Delta_j(10|00)\} + D_j(11)\Delta_j(10|00) + D_{j'}(11)\Delta_j(01|00) + Y_j(00)]
\end{aligned}$$



Doing what has been done with the three previous rows:

$$\begin{aligned}
\mathbb{E} [Z_j Z_{j'} Y_j] &= P^2 \{ \mathbb{E} [\Delta_j(11|01) - \Delta_j(10|00) \mid D_j(11) = 1, D_{j'}(11) = 1] \mathbb{P}(D_j(11) = 1, D_{j'}(11) = 1) \\
&\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1] \mathbb{P}(D_j(11) = 1) \\
&\quad + \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 0] \mathbb{P}(D_{j'}(11) = 1) \} \\
&\quad + P^2 \mathbb{E} [Y_j(00)].
\end{aligned}$$

Before  $\mathbf{M}^{-1}\mathbf{v}$  is computed in Appendix C.5, I introduce new notation for further simplification:

$$\begin{aligned}
e_1 &\equiv \mathbb{E} [\Delta_j(11|01) - \Delta_j(10|00) \mid D_j(11) = 1, D_{j'}(11) = 1] \\
e_2 &\equiv \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \\
e_3 &\equiv \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 0] \\
e_4 &\equiv \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1] \\
e_5 &\equiv \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(01) = 1] \\
e_6 &\equiv \mathbb{E} [Y_j(00)] \\
e_7 &\equiv \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1] \\
e_8 &\equiv \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1].
\end{aligned}$$

Recalling the notation in Appendix C.3,

$$\begin{aligned}
q &\equiv \mathbb{P}(D_j(11) = 1) \\
\lambda q &\equiv \mathbb{P}(D_{j'}(11) = 1) \\
r &\equiv \mathbb{P}(D_j(11) = 1, D_j(10) = 0) \\
\mu r &\equiv \mathbb{P}(D_{j'}(11) = 1, D_{j'}(01) = 0) \\
\bar{q} &\equiv \mathbb{P}(D_j(11) = 1, D_{j'}(11) = 1) \\
q - r &= \mathbb{P}(D_j(10) = 1) = \mathbb{P}(D_j(10) = 1, D_j(11) = 1) \quad (\text{last equality by (16)}) \\
\lambda q - \mu r &= \mathbb{P}(D_{j'}(01) = 1) = \mathbb{P}(D_{j'}(01) = 1, D_{j'}(11) = 1) \quad (\text{last equality by (17)})
\end{aligned}$$

we can express  $\mathbf{v}$  as

$$\mathbf{v} = \begin{bmatrix} \mathbb{E}[Y_j] \\ \mathbb{E}[Z_j Y_j] \\ \mathbb{E}[Z_{j'} Y_j] \\ \mathbb{E}[Z_j Z_{j'} Y_j] \end{bmatrix} = \begin{bmatrix} P^2\{e_1 \bar{q} + e_2 r + e_3 \mu r\} + P e_4 (q - r) + P e_5 (\lambda q - \mu r) + e_6 \\ P^2\{e_1 \bar{q} + e_2 r + e_7 \lambda q\} + P e_4 (q - r) + P e_6 \\ P^2\{e_1 \bar{q} + e_8 q + e_3 \mu r\} + P e_5 (\lambda q - \mu r) + P e_6 \\ P^2\{e_1 \bar{q} + e_8 q + e_7 \lambda q\} + P^2 e_6 \end{bmatrix}. \quad (24)$$

### C.5. Computing $M^{-1}\mathbf{v}$

Having computed  $M$  and  $\mathbf{v}$  in Appendix C.3 and C.4, the final step in the proof of Theorem 2 is evaluating  $M^{-1}\mathbf{v} \in \mathbb{R}^{4 \times 1}$ . In this 4-long column vector, each row is a linear combination of the elements of  $\mathbf{v}$ , where the weights are the elements in the corresponding rows of  $M^{-1}$ . A6 Invertibility ensures that  $M^{-1}$  exists (see the end of C.3). Suppose that

$$M^{-1} = \begin{bmatrix} m_1^{(1)} & m_2^{(1)} & m_3^{(1)} & m_4^{(1)} \\ m_1^{(2)} & m_2^{(2)} & m_3^{(2)} & m_4^{(2)} \\ m_1^{(3)} & m_2^{(3)} & m_3^{(3)} & m_4^{(3)} \\ m_1^{(4)} & m_2^{(4)} & m_3^{(4)} & m_4^{(4)} \end{bmatrix},$$

then the  $r$ th row for  $r = 1, 2, 3, 4$  in  $M^{-1}\mathbf{v}$ ,  $(M^{-1}\mathbf{v})_r$  is given by

$$(M^{-1}\mathbf{v})_r = m_1^{(r)}v_1 + m_2^{(r)}v_2 + m_3^{(r)}v_3 + m_4^{(r)}v_4.$$

If we substitute in the results from (24) we obtain

$$\begin{aligned} (M^{-1}\mathbf{v})_r &= m_1^{(r)}P^2\{e_1\bar{q} + e_2r + e_3\mu r\} + Pe_4(q - r) + Pe_5(\lambda q - \mu r) + e_6 \\ &+ m_2^{(r)}[P^2\{e_1\bar{q} + e_2r + e_7\lambda q\} + Pe_4(q - r) + Pe_6] \\ &+ m_3^{(r)}[P^2\{e_1\bar{q} + e_8q + e_3\mu r\} + Pe_5(\lambda q - \mu r) + Pe_6] \\ &+ m_4^{(r)}[P^2\{e_1\bar{q} + e_8q + e_7\lambda q\} + P^2e_6], \end{aligned}$$

and after collecting terms

$$\begin{aligned} (M^{-1}\mathbf{v})_r &= e_1[P^2\bar{q}\{m_1^{(r)} + m_2^{(r)} + m_3^{(r)} + m_4^{(r)}\}] \\ &+ e_2[P^2r(m_1^{(r)} + m_2^{(r)})] \\ &+ e_3[P^2\mu r(m_1^{(r)} + m_3^{(r)})] \\ &+ e_4[P(q - r)(m_1^{(r)} + m_2^{(r)})] \\ &+ e_5[P(\lambda q - \mu r)(m_1^{(r)} + m_3^{(r)})] \\ &+ e_6[m_1^{(r)} + Pm_2^{(r)} + Pm_3^{(r)} + P^2m_4^{(r)}] \\ &+ e_7[P^2\lambda q(m_2^{(r)} + m_4^{(r)})] \\ &+ e_8[P^2q(m_3^{(r)} + m_4^{(r)})] \\ &\equiv \sum_{k=1}^8 e_k w_k^{(r)}. \end{aligned}$$

At this point, to compute  $w$ 's, it is convenient to rely on a symbolic math package, SymPy. The code in `symbolic_inversion.py` computes  $\{w_k^{(r)}\}_{k=1}^8$  for  $r = 1, 2, 3, 4$ , and the values are shown in Table 2.

**Table 2: Output from `symbolic_inversion.py`**

$\{w_k^{(r)}\}_{k=1}^8 \setminus \text{Row}$	1	2	3	4
$w_1^{(r)}$	0	0	0	1
$w_2^{(r)}$	0	$\frac{Pr}{(q-r)(1-P)}$	0	$\frac{-Pqr}{\bar{q}(q-r)(1-P)}$
$w_3^{(r)}$	0	0	$\frac{P\mu r}{(\lambda q - \mu r)(1-P)}$	$\frac{-P\lambda q\mu r}{\bar{q}(\lambda q - \mu r)(1-P)}$
$w_4^{(r)}$	0	$\frac{q-r}{(q-r)(1-P)}$	0	$\frac{-q(q-r)}{\bar{q}(q-r)(1-P)}$
$w_5^{(r)}$	0	0	$\frac{\lambda q - \mu r}{(\lambda q - \mu r)(1-P)}$	$\frac{-\lambda q(\lambda q - \mu r)}{\bar{q}(\lambda q - \mu r)(1-P)}$
$w_6^{(r)}$	1	0	0	0
$w_7^{(r)}$	0	0	$\frac{-P\lambda q}{(\lambda q - \mu r)(1-P)}$	$\frac{\lambda q(P\mu r + \lambda q - \mu r)}{\bar{q}(\lambda q - \mu r)(1-P)}$
$w_8^{(r)}$	0	$\frac{-Pq}{(q-r)(1-P)}$	0	$\frac{q(Pr + q - r)}{\bar{q}(q-r)(1-P)}$

In the following, I plug in these values to  $\sum_{k=1}^8 e_k w_k^{(r)}$  to evaluate  $\mathbf{M}^{-1}\mathbf{v}$  row by row (weights taking on the value 0 are not written out). First, examine expressions in the denominator.

## Denominator

By Lemma 6 and  $P \in (0, 1)$ , the denominator is nonzero. Furthermore,

$$q - r = \mathbb{P}(D_j(10) = 1, D_j(11) = 1) \quad (25)$$

$$\lambda q - \mu r = \mathbb{P}(D_{j'}(01) = 1, D_{j'}(11) = 1) \quad (26)$$

by (16) and (17).

**Row of  $M^{-1}\mathbf{v}$ : 1st,  $(M^{-1}\mathbf{v})_1$ .**

$$(M^{-1}\mathbf{v})_1 = w_6^{(1)}e_6 = 1e_6 = \mathbb{E}[Y_j(00)].$$

Hence the first claim of Theorem 2 is proved.

**Row of  $M^{-1}\mathbf{v}$ : 2nd,  $(M^{-1}\mathbf{v})_2$ .**

$$\begin{aligned} (M^{-1}\mathbf{v})_2 &= w_2^{(2)}e_2 + w_4^{(2)}e_4 + w_8^{(2)}e_8 \\ &= \frac{1}{(q-r)(1-P)} \{Pr \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] + (q-r) \mathbb{E}[\Delta_j(10|00) \mid D_j(10) = 1] - Pq \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1]\} \\ &= \frac{1}{(q-r)(1-P)} \{P[r \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] - q \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1]] + (q-r) \mathbb{E}[\Delta_j(10|00) \mid D_j(10) = 1]\} \end{aligned} \quad (27)$$

To figure out what  $q \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1]$  is, expand  $\mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1]$  using the law of total probability and the definition of  $q = \mathbb{P}(D_j(11) = 1)$  and conditional probability.

$$\begin{aligned}
\mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1] &= \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \mathbb{P}(D_j(10) = 0 \mid D_j(11) = 1) \\
&\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] \mathbb{P}(D_j(10) = 1 \mid D_j(11) = 1) \\
&= \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \frac{\mathbb{P}(D_j(10) = 0, D_j(11) = 1)}{\mathbb{P}(D_j(11) = 1)} \\
&\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] \frac{\mathbb{P}(D_j(10) = 1, D_j(11) = 1)}{\mathbb{P}(D_j(11) = 1)} \\
&= \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \frac{\mathbb{P}(D_j(10) = 0, D_j(11) = 1)}{q} \\
&\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] \frac{\mathbb{P}(D_j(10) = 1, D_j(11) = 1)}{q}
\end{aligned}$$

45

which in turn implies

$$\begin{aligned}
q \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1] &= \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \mathbb{P}(D_j(10) = 0, D_j(11) = 1) \\
&\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] \mathbb{P}(D_j(10) = 1, D_j(11) = 1). \tag{28}
\end{aligned}$$

Applying the same expansion to  $(q - r) \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1]$ , noting that  $q - r = \mathbb{P}(D_j(10) = 1)$  (see (14)) leads to

$$\begin{aligned}
(q - r) \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1] &= \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 1] \mathbb{P}(D_j(11) = 1, D_j(10) = 1) \\
&\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 0] \mathbb{P}(D_j(11) = 0, D_j(10) = 1)
\end{aligned}$$

where by A5 Monotonicity  $\mathbb{P}(D_j(11) = 0, D_j(10) = 1) = 0$  so

$$(q - r) \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1] = \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 1] \mathbb{P}(D_j(11) = 1, D_j(10) = 1) \quad (29)$$

Plugging in the above expressions, the definition of  $r$ , and the denominator, into (27):

$$\begin{aligned} (\mathbf{M}^{-1}\mathbf{v})_2 &= \frac{1}{(q-r)(1-P)} \{P[r \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] - q \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1]] + (q-r) \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1]\} \\ &= \frac{1}{(q-r)(1-P)} \{P\{\mathbb{P}(D_j(11) = 1, D_j(10) = 0) \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \\ &\quad - [\mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \mathbb{P}(D_j(10) = 0, D_j(11) = 1) \\ &\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] \mathbb{P}(D_j(10) = 1, D_j(11) = 1)]\} \\ &\quad + \mathbb{P}(D_j(11) = 1, D_j(10) = 1) \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 1]\} \\ &= \frac{1}{(q-r)(1-P)} (1-P) \mathbb{P}(D_j(11) = 1, D_j(10) = 1) \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 1] \end{aligned}$$

Finally, substituting in  $q - r = \mathbb{P}(D_j(11) = 1, D_j(10) = 1)$  (see (25)) gives us

$$\begin{aligned} (\mathbf{M}^{-1}\mathbf{v})_2 &= \frac{1}{\mathbb{P}(D_j(11) = 1, D_j(10) = 1) (1-P)} (1-P) \mathbb{P}(D_j(11) = 1, D_j(10) = 1) \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 1] \\ &= \mathbb{E} [\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 1]. \end{aligned}$$

Hence the second claim of Theorem 2 is proved.

**Row of  $\mathbf{M}^{-1}\mathbf{v}$ : 3rd,  $(\mathbf{M}^{-1}\mathbf{v})_3$ .**

Exploiting symmetry and that  $\lambda q - \mu r = \mathbb{P}(D_{j'}(01) = 1, D_{j'}(11) = 1)$  (see (26)), we can proceed in the same way as with the second row.

$$\begin{aligned}
(\mathbf{M}^{-1}\mathbf{v})_3 &= w_3^{(3)}e_3 + w_5^{(3)}e_5 + w_7^{(3)}e_7 \\
&= \frac{1}{(\lambda q - \mu r)(1 - P)} \{ P\mu r \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 0] + (\lambda q - \mu r) \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(01) = 1] - P\lambda q \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1] \} \\
&= \frac{1}{(\lambda q - \mu r)(1 - P)} \{ P[\mu r \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 0] - \lambda q \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1]] + (\lambda q - \mu r) \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(01) = 1] \} \\
&= \frac{1}{(\lambda q - \mu r)(1 - P)} (1 - P) \mathbb{P} (D_{j'}(11) = 1, D_{j'}(01) = 1) \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(01) = 1, D_{j'}(11) = 1] \\
&= \frac{1}{(1 - P) \mathbb{P} (D_{j'}(11) = 1, D_{j'}(01) = 1)} (1 - P) \mathbb{P} (D_{j'}(11) = 1, D_{j'}(01) = 1) \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(01) = 1, D_{j'}(11) = 1] \\
&= \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(01) = 1, D_{j'}(11) = 1].
\end{aligned}$$

Hence the third claim of Theorem 2 is proved.

**Row of  $\mathbf{M}^{-1}\mathbf{v}$ : 4th,  $(\mathbf{M}^{-1}\mathbf{v})_4$ .**

For the last element we have the most non-zero weights in Table 2, amounting to

$$\begin{aligned}
(\mathbf{M}^{-1}\mathbf{v})_4 &= w_1^{(4)}e_1 + w_2^{(4)}e_2 + w_3^{(4)}e_3 + w_4^{(4)}e_4 + w_5^{(4)}e_5 + w_7^{(4)}e_7 + w_8^{(4)}e_8 \\
&= w_1^{(4)}e_1 + (w_2^{(4)}e_2 + w_4^{(4)}e_4 + w_8^{(4)}e_8) + (w_3^{(4)}e_3 + w_5^{(4)}e_5 + w_7^{(4)}e_7).
\end{aligned}$$

Grouping the sum as indicated with the brackets renders the computation easier. In the following, I evaluate the sum group by group.

- $w_1^{(4)}e_1$

$$w_1^{(4)}e_1 = 1e_1 = \mathbb{E} [\Delta_j(11|01) - \Delta_j(10|00) \mid D_j(11) = 1, D_{j'}(11) = 1] \quad (30)$$



- $w_2^{(4)}e_2 + w_4^{(4)}e_4 + w_8^{(4)}e_8$

$$\begin{aligned} w_2^{(4)}e_2 + w_4^{(4)}e_4 + w_8^{(4)}e_8 &= \frac{1}{\bar{q}(q-r)(1-P)} \{-Pqre_2 - (q-r)qe_4 + q(Pr + q-r)e_8\} \\ &= \frac{1}{\bar{q}(q-r)(1-P)} \{Pqr(e_8 - e_2) + (q-r)q(e_8 - e_4)\}. \end{aligned}$$

Next, I compute  $Pqr(e_8 - e_2)$  and  $(q-r)q(e_8 - e_4)$ .

- $Pqr(e_8 - e_2)$  From the definition of  $e_8$  and from (28), we know that

$$\begin{aligned} qre_8 &= qr \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1] \\ &= \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \mathbb{P}(D_j(10) = 0, D_j(11) = 1) r \\ &\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] \mathbb{P}(D_j(10) = 1, D_j(11) = 1) r, \end{aligned}$$

where by definition  $\mathbb{P}(D_j(10) = 0, D_j(11) = 1) = r$ , and by (25)  $\mathbb{P}(D_j(11) = 1, D_j(10) = 1) = q - r$  so

$$\begin{aligned} qr \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1] &= \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \mathbb{P}(D_j(10) = 0, D_j(11) = 1) r \\ &\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] \mathbb{P}(D_j(10) = 1, D_j(11) = 1) r \\ &= \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] r^2 \\ &\quad + \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] (q-r)r. \end{aligned}$$

By the definition of  $e_2$

$$qre_2 = \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] qr$$

so

$$\begin{aligned}
Pqr(e_8 - e_2) &= P\{\mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] r^2 \\
&\quad + \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] (q - r)r \\
&\quad - \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] qr\} \\
&= P\{\mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] (q - r)r - \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] (q - r)r\} \\
&= P(q - r)r\{\mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] - \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0]\}.
\end{aligned}$$

••  $(q - r)q(e_8 - e_4)$  From the definition of  $e_8$  and from (28), we know that

$$\begin{aligned}
(q - r)qe_8 &= (q - r)q \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1] \\
&= \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] \mathbb{P}(D_j(10) = 0, D_j(11) = 1) (q - r) \\
&\quad + \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] \mathbb{P}(D_j(10) = 1, D_j(11) = 1) (q - r)
\end{aligned}$$

and from (25)  $\mathbb{P}(D_j(10) = 1, D_j(11) = 1) = q - r$ , and by definition  $\mathbb{P}(D_j(10) = 0, D_j(11) = 1) = r$  so

$$(q - r)qe_8 = \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] r(q - r) + \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] (q - r)^2.$$

From the definition of  $e_4$  and (29)

$$\begin{aligned}
(q - r)qe_4 &= (q - r)q \mathbb{E}[\Delta_j(10|00) \mid D_j(10) = 1] \\
&= \mathbb{E}[\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 1] \mathbb{P}(D_j(11) = 1, D_j(10) = 1) q
\end{aligned}$$

then by (25)  $\mathbb{P}(D_j(10) = 1, D_j(11) = 1) = q - r$  so

$$(q - r)qe_4 = \mathbb{E}[\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 1](q - r)q$$

It follows that

$$\begin{aligned} (q - r)q(e_8 - e_4) &= \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0]r(q - r) + \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1](q - r)^2 \\ &\quad - q(q - r)\mathbb{E}[\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 1] \\ &= \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1](q - r)(q - r - q) + \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0]r(q - r) \\ &= r(q - r)\{\mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] - \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1]\} \end{aligned}$$

Having computed  $Pqr(e_8 - e_2)$  and  $(q - r)q(e_8 - e_4)$  we can obtain

$$\begin{aligned} Pqr(e_8 - e_2) + (q - r)q(e_8 - e_4) &= P(q - r)r\{\mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1] - \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0]\} \\ &\quad + r(q - r)\{\mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] - \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1]\} \\ &= (q - r)r(1 - P)\{\mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] - \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1]\} \end{aligned}$$

Thus

$$\begin{aligned} w_2^{(4)}e_2 + w_4^{(4)}e_4 + w_8^{(4)}e_8 &= \frac{1}{\bar{q}(q - r)(1 - P)}\{Pqr(e_8 - e_2) + (q - r)q(e_8 - e_4)\} \\ &= \frac{1}{\bar{q}(q - r)(1 - P)}(q - r)r(1 - P)\{\mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] - \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1]\} \\ &= \frac{r}{\bar{q}}\{\mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] - \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1]\} \end{aligned}$$

- $w_3^{(4)}e_3 + w_5^{(4)}e_5 + w_7^{(4)}e_7$

$$\begin{aligned} w_3^{(4)}e_3 + w_5^{(4)}e_5 + w_7^{(4)}e_7 &= \frac{1}{\bar{q}(\lambda q - \mu r)(1 - P)} \{-P\lambda q \mu r e_3 - \lambda q(\lambda q - \mu r)e_5 + \lambda q(P\mu r + \lambda q - \mu r)e_7\} \\ &= \frac{1}{\bar{q}(\lambda q - \mu r)(1 - P)} \{P\lambda q \mu r(e_7 - e_3) + \lambda q(\lambda q - \mu r)(e_7 - e_5)\}. \end{aligned}$$

Due to the symmetry of the problem and the definitions of  $e$ 's,  $\lambda q, \mu r$  this sum behaves exactly as  $Pqr(e_8 - e_2) + (q - r)q(e_8 - e_4)$ , only that this time we have  $D_{j'}(11)$  instead of  $D_j(11)$ ,  $D_{j'}(01)$  instead of  $D_j(10)$  and  $\Delta_j(01|00)$  instead of  $\Delta_j(10|00)$ . So the proof is exactly the same, and we have

$$w_3^{(4)}e_3 + w_5^{(4)}e_5 + w_7^{(4)}e_7 = \frac{\mu r}{\bar{q}} \{\mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 0] - \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 1]\}.$$

51

Putting all this together leaves us with

$$\begin{aligned} (\mathbf{M}^{-1}\mathbf{v})_4 &= w_1^{(4)}e_1 + (w_2^{(4)}e_2 + w_4^{(4)}e_4 + w_8^{(4)}e_8) + (w_3^{(4)}e_3 + w_5^{(4)}e_5 + w_7^{(4)}e_7) \\ &= \mathbb{E} [\Delta_j(11|01) - \Delta_j(10|00) \mid D_j(11) = 1, D_{j'}(11) = 1] \\ &\quad + \frac{r}{\bar{q}} \{\mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] - \mathbb{E} [\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1]\} \\ &\quad + \frac{\mu r}{\bar{q}} \{\mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 0] - \mathbb{E} [\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 1]\} \end{aligned}$$

Hence the fourth statement of Theorem 2 is proved.

## C.6. Overview of the Proof

Let

$$\begin{aligned}\hat{\boldsymbol{\theta}}_j &\equiv (n^{-1} \mathbf{Z}'_j \mathbf{D}_j)^{-1} n^{-1} \mathbf{Z}'_j \mathbf{y}_j \\ &= \left( n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} \mathbf{d}'_{ji} \right)^{-1} n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} y_{ji} \\ \mathbf{d}_{ji} &\equiv [1, d_{ji}, d_{j'i}, d_{ji} d_{j'i}]' \in \{0, 1\}^{4 \times 1} \\ \mathbf{z}_{ji} &\equiv [1, z_{ji}, z_{j'i}, z_{ji} z_{j'i}]' \in \{0, 1\}^{4 \times 1}\end{aligned}$$

for  $j \in \{A, B\}$  and  $j' \in \{A, B\} \setminus \{j\}$ .

Examining  $\text{plim } \hat{\boldsymbol{\theta}}_j$ , i.e. the probability limit of  $\hat{\boldsymbol{\theta}}_j$ , or equivalently  $\boldsymbol{\theta}_j : \lim_{n \rightarrow \infty} \mathbb{P} \left( \|\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j\|_2^2 > \varepsilon \right) = 0$  for any  $\varepsilon > 0$  leads to

$$\begin{aligned}\text{plim } \hat{\boldsymbol{\theta}}_j &= \text{plim} \left( n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} \mathbf{d}'_{ji} \right)^{-1} n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} y_{ji} \\ &= \left( \text{plim } n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} \mathbf{d}'_{ji} \right)^{-1} \text{plim } n^{-1} \sum_{i=1}^n \mathbf{z}_{ji} y_{ji} \\ &= \mathbb{E} \left[ \mathbf{z}_{ji} \mathbf{d}'_{ji} \right]^{-1} \mathbb{E} \left[ \mathbf{z}_{ji} y_{ji} \right] \\ &= \mathbb{E} \left[ \begin{bmatrix} 1 \\ Z_j \\ Z_{j'} \\ Z_j Z_{j'} \end{bmatrix} \begin{bmatrix} 1 & D_j & D_{j'} & D_j D_{j'} \end{bmatrix} \right]^{-1} \mathbb{E} \left[ \begin{bmatrix} 1 \\ Z_j \\ Z_{j'} \\ Z_j Z_{j'} \end{bmatrix} Y_j \right] \\ &\equiv \mathbf{M}^{-1} \mathbf{v}\end{aligned}$$

by the continuous mapping property of the probability limit and by the Weak Law of Large Numbers for i.i.d. data (as the data across pairs are i.i.d.).

I proved that Identifying Assumptions  $A1 - A6$  are sufficient to establish

$$\begin{aligned}
(\mathbf{M}^{-1}\mathbf{v})_1 &= \mathbb{E}[Y_j(00)] \\
(\mathbf{M}^{-1}\mathbf{v})_2 &= \mathbb{E}[\Delta_j(10|00) \mid D_j(10) = 1, D_j(11) = 1] \\
(\mathbf{M}^{-1}\mathbf{v})_3 &= \mathbb{E}[\Delta_j(01|00) \mid D_{j'}(01) = 1, D_{j'}(11) = 1] \\
(\mathbf{M}^{-1}\mathbf{v})_4 &= \mathbb{E}[\Delta_j(11|01) - \Delta_j(10|00) \mid D_j(11) = 1, D_{j'}(11) = 1] \\
&\quad + \frac{r}{\bar{q}}\{\mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 0] - \mathbb{E}[\Delta_j(10|00) \mid D_j(11) = 1, D_j(10) = 1]\} \\
&\quad + \frac{\mu r}{\bar{q}}\{\mathbb{E}[\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 0] - \mathbb{E}[\Delta_j(01|00) \mid D_{j'}(11) = 1, D_{j'}(01) = 1]\}.
\end{aligned}$$

53

Thus the proof of Theorem 2 is complete ■