

# Robust Testing for Causal Inference in Observational Studies

Md. Noor-E-Alam

Dept. of Mechanical and Industrial Engineering, Northeastern University, 360 Huntington Avenue, Boston, MA 02115, USA  
md.alam@neu.edu

Cynthia Rudin

MIT CSAIL and Sloan School of Management, Massachusetts Institute of Technology, Cambridge, MA 02139, USA  
rudin@mit.edu

A vast number of causal inference studies use matching techniques, where treatment cases are matched with similar control cases. For observational data in particular, we claim there is a major source of uncertainty that is essentially ignored in these tests, which is the way the assignments of matched pairs are constructed. It is entirely possible, for instance, that a study reporting an estimated treatment effect with  $P$ -value of  $10^{-4}$  can be redone in almost the same way, with the same match quality, yet with a  $P$ -value well above 0.10, making the test result no longer significant. Experimenters often specifically choose *not* to consider the output in the design of the assignments; this allows for easier computation and clearer testing, but it does not consider possible biases in the way the assignments were constructed. What we would really like to be able to report is that *no matter* which assignment we choose, as long as the match is sufficiently good, then the hypothesis test result still holds. This will be called a *robust* matched pairs test, since its result is robust to the choice of the assignment. In this paper, we provide methodology based on discrete optimization to create these robust tests. This method explores the full variation of possible results one can obtain with all possible acceptable assignments. It demonstrates that one cannot always trust statistically significant results, even with a large number of matched pairs.

*Key words*: causal inference, observational studies, hypothesis test, matched pairs design, discrete optimization, integer programming.

---

## 1. Introduction

As massive and varied amounts of observational data are accumulating in healthcare, internet marketing, and governance, these data are increasingly used for understanding important cause and effect relationships. We might want to know whether a policy causes people to use fewer public services, or we might want to know whether a particular drug causes a side effect, or whether a view of an internet advertisement results in an increased chance of purchase. Controlled trials with matching treatment and controlled populations are often small and expensive, and not often possible due to the ethics of offering the treatment (e.g., exposure to lead), or perhaps not possible due to the fact that the treatment happened only in the past (e.g., lived in Morristown in the 1960s). Typically, for large amounts of observational data, we can construct matched pairs of observations, where one received the treatment and the other did not. Observations are matched on a set of attributes (age, gender, race), propensity scores, etc. The quality of the matches can be very important – poor matching algorithms can lead to wrong conclusions, and potentially to harmful politics, laws, and medical decisions.

Classically, assignments of treatment and control units to matches are constructed using a fixed design (Rosenbaum 2010a), without regard to the outcome (Rubin 2007, 2008). From our perspective, this could be a major flaw in the current paradigm. Choosing a single fixed design ignores a major source of uncertainty, which is the design itself, or in other words, the uncertainty related to the choice of experimenter. What if there were two possible equally good matchings, one where the treatment effect estimate is very strong and one where it is nonexistent? When we report a result on a particular matching assignment, we thus ignore the possibility of the opposite result occurring on an equally good assignment. It is entirely possible that two separate researchers studying the same effect on the same data, using two different equally good sets of pairs, would get results that disagree.

Our goal is to create robust matched pairs hypothesis tests for observational data. These tests implicitly consider *all possible reasonably good assignments* and consider the *range of possible*

*outcomes* for tests on these data. This is a more computationally demanding approach to hypothesis testing than the standard approach where one considers just a single assignment, but the result would be robust to the choice of experimenter. It is computationally infeasible (and perhaps not very enlightening) to explicitly compute all possible assignments, but it is possible to look at the range of outcomes associated with them. In particular, our algorithms compute the maximum and minimum of quantities like the treatment effect estimate,  $P$ -value, and  $z$ -score among the set of reasonably good assignments.

Finding a set of matched pairs that obey certain conditions is purely a data mining problem – we aim to locate a particular pattern of data within a database. Similar subfields of data mining, where the goal is to locate optimal subsets of data, include association rule mining and event detection. For these types of problems, modern discrete optimization techniques have rarely been used, though there is one recent precedent in the literature for matched pairs, namely the line of work by Zubizarreta (2012), Zubizarreta et al. (2013, 2014). Optimization techniques have major advantages over other types of approaches: (i) they are extremely flexible and allow the experimenter to match on very complex conditions, such as quantiles of attributes, which network optimization methods for matching cannot handle; (ii) they can be computationally efficient, depending on the strength of the integer programming formulation – strong formulations have relaxations close to the set of feasible integer solutions, and (iii) mixed-integer linear programs have guarantees on the optimality of the solution – in fact, they produce upper and lower bounds on the value of the optimal solution.

In what follows, we provide three basic constructions for robust hypothesis tests, one for a matched pairs  $z$ -test on the difference between two means, the second one for a McNemar’s test for proportion data (1:1 matched pair design), and the third one for a  $\chi^2$  test also for proportion data (1:m matched pair design). The  $z$ -test and  $\chi^2$  test create nonlinear optimization problems that (in our experiments) were not able to be solved with a MINLP (mixed integer nonlinear programming) solver within a reasonable amount of time. Instead, we propose algorithms that solve a series of integer linear programs (ILPs) along a coarse one-dimensional mesh and then at finer

one-dimensional meshes, until we achieve a solution with the desired precision. Computationally these methods are much more attractive, and scale to large datasets, as we demonstrate.

The remainder of this paper is organized as follows. Section 2 formalizes robust testing as our framework defines it. In Section 3, we propose an optimization-based approach to carry out the robust  $z$ -test for real valued outcome data. In Sections 4 and 5, we propose two optimization-based approaches for inference on a proportion, one for 1:1 pair matching and another one is for 1-to-many matching. Section 6 presents empirical results and analysis.

## 2. Matching for Robust Tests

This work concerns the potential outcomes framework (see Holland 1986, Rubin 1974). In our notation  $X$  represents a vector of covariates for an individual, and  $Y$  is an outcome variable that depends on  $X$  and whether the patient was treated,  $Y(1, X)$  is the random variable for the outcome of a treated individual with covariates  $X$ , and  $Y(0, X)$  is the outcome for an untreated individual. We make the classical SUTVA assumption (the treatment status of any unit does not affect the potential outcomes of the other units), and assume conditional ignorability, which is that  $Y$  is independent of treatment  $T$  given  $X$ . We also assume unconfoundedness. Our goal is to determine whether we can reject the claim that certain standard quantities are zero, such as:

$$\text{ATT: } \mathbb{E}_{(X|T=1), (Y|T=1, X)}[Y(1, X) - Y(0, X)|T = 1]$$

$$\text{ATE: } \mathbb{E}_{(X, Y)}[Y(1, X) - Y(0, X)].$$

The distribution of  $X$  is different for the treatment and control groups, because there is a bias as to who receives the treatment. The distribution of  $Y|0, X$  and  $Y|1, X$  will be different if the treatment has an effect. To conduct our hypothesis test we have observations:

$$(\mathbf{x}_1^t, y_1^t), \dots, (\mathbf{x}_T^t, y_T^t)$$

$$(\mathbf{x}_1^c, y_1^c), \dots, (\mathbf{x}_C^c, y_C^c).$$

A matching operator determines which control is assigned to which treatment. For simplicity, we write an assignment operator that assigns at most one control to each treatment without

replacement, though we loosen this definition later.

*Definition:* A matching operator  $\Omega: \{1, \dots, T\} \rightarrow \{1, \dots, C, \emptyset\}$  obeys the following: if  $i \neq k$  and  $\Omega(i) \neq \emptyset$  then  $\Omega(i) \neq \Omega(k)$ . That is, no two treatment units  $i$  and  $k$  are assigned to the same control unit. We define the size of the matching to be  $|\Omega| = \sum_i \mathbf{1}_{[\Omega(i) \neq \emptyset]}$ . The set of all matching operators is  $\mathcal{A}$  = set of all assignments  $\{\Omega\}$ .

## 2.1. Traditional and Current Matching Procedures

Traditional and current matching procedures perform the following two steps, where the matching procedure (denoted “Algorithm” or “Alg” below) is usually not denoted explicitly. (For what follows, we need to denote this explicitly.) The test statistic is a function of the matching procedure.

### Classical Procedure.

Compute match assignment :  $\Omega_{\text{Alg}} = \text{Algorithm}(\{\mathbf{x}_i^t\}_{i=1}^T, \{\mathbf{x}_i^c\}_{i=1}^C)$

Compute causal effect and test statistic :  $\hat{E}[\bar{y}^t - \bar{y}^c], z_{\Omega_{\text{Alg}}}(\{y_i^t\}_{i=1}^T, \{y_i^c\}_{i=1}^C)$ , and  $\text{P-value}(z_{\Omega_{\text{Alg}}})$ .

Here we gave an example of a matched pairs test for the difference in means. Examples of the matching part of the Classical Procedure could include the following:

#### Matching Example 1:

$$\Omega_{\text{Alg}} = \text{Greedy Matching Algorithm}(\{(\mathbf{x}_i^t)\}_{i=1}^T, \{(\mathbf{x}_i^c)\}_{i=1}^C),$$

where the first treatment observation  $\mathbf{x}_1^t$  is matched to its nearest control, then the second treatment observation is matched, and so on. In this case, the test statistic and P-value for the ATT problem would depend on the greedy procedure. This means the P-value would be very sensitive to the particular experimenter and setup – it would even depend on the order in which the data were recorded.

**Matching Example 2:** Let us consider the classical personnel assignment problem, which is approximately or exactly solved by several software packages:

$$\begin{aligned} \Omega_{\text{Alg}} &\in \text{Optimal Assignment Matching Algorithm}(\{(\mathbf{x}_i^t)\}_{i=1}^T, \{(\mathbf{x}_i^c)\}_{i=1}^C), \\ &= \underset{\Omega \in \mathcal{A}}{\text{argmin}} \left( \sum_{i=1}^T \text{dist}(\mathbf{x}_i^t, \mathbf{x}_{\Omega(i)}^c) \right), \end{aligned}$$

where every treatment case must be matched to one or more controls for the ATT problem. (That is,  $\text{dist}(\mathbf{x}_i^t, \mathbf{x}_\emptyset^c) = \infty$ .) In this case the test statistic and P-value depend on the choice of distance measure and the result of this optimization procedure. If there are multiple optimal solutions to the matching problem, or close-to-optimal solutions, they would not be considered in the P-value calculation, and thus the result could depend upon the experimenter and/or the computer memory layout or specific instantiation of the algorithm. It could also depend heavily on outliers that are difficult to match. See Rosenbaum (1989) for a review of the assignment problem.

**Matching Example 3:** Let us consider a more flexible matching procedure, which is a simple version of Problem 1 of Rosenbaum (2012). Rosenbaum’s work points out that some treated subjects may be too extreme to match. His matching algorithm makes three optimal decisions at once: (i) the number of treated subjects to match, (ii) the identity of the treated subjects to match, and (iii) the identity of the controls with whom they are paired.

$$\Omega_{\text{Alg}} \in \operatorname{argmin}_{\Omega \in \mathcal{A}} \left( \sum_{i=1}^{|\Omega|} \text{dist}(\mathbf{x}_i^t, \mathbf{x}_{\Omega(i)}^c) : |\Omega| \geq N \right),$$

where  $|\Omega|$  is the number of matches, which we want to be sufficiently large. This aims at a distribution of “marginal patients” for which there is an overlap in density between treatment and control patients, rather than either the density of the ATT or ATE problems. Rosenbaum (2012) also discusses the importance of covariate balance. In our formulations, we assume additional constraints for covariate balance would be included. A different approach is taken by Ho et al. (2007) who remove points outside of the overlap as a pre-processing step. In this work one could use Ho et al. (2007) as a pre-processing step, but our formulations encompass the full setting of Rosenbaum (2012).

Again in this example, a single assignment is chosen, and no assignments that are slightly suboptimal are considered. It does not consider, for instance, all solutions for which the distance measure used in the procedure is slightly different, all solutions for which the solution is slightly suboptimal, or all solutions for which  $N$  is varied within a reasonable range. There could be thousands or millions of assignments that are slightly different than the one chosen by this procedure, all

leading to possibly different P-values. If all of these P-values were below (for instance) 0.05, the result would be more robust than if only some of them were.

Before we move on, let us give an example of the test statistic computation with the extra notation to make the matching explicit.

**Test Statistic Computation Example:** We consider the 2-sample matched pair  $z$ -test.

$$\begin{aligned}\bar{d}_{\Omega_{\text{Alg}}} &= \frac{1}{|\Omega_{\text{Alg}}|} \sum_{i=1}^{|\Omega_{\text{Alg}}|} (y_i^t - y_{\Omega_{\text{Alg}}(i)}^c) \\ \sigma_{\Omega_{\text{Alg}}}^2 &\approx s_{\Omega_{\text{Alg}}}^2 = \frac{1}{|\Omega_{\text{Alg}}| - 1} \sum_{i=1}^{|\Omega_{\text{Alg}}|} (y_i^t - y_{\Omega_{\text{Alg}}(i)}^c - \bar{d}_{\Omega_{\text{Alg}}})^2 \\ z_{\Omega_{\text{Alg}}} &= \frac{\bar{d}_{\Omega_{\text{Alg}}} \sqrt{n}}{\sigma_{\Omega_{\text{Alg}}}}, \quad \text{P-value}_{\Omega_{\text{Alg}}} = 1 - \Phi(z_{\Omega_{\text{Alg}}}).\end{aligned}$$

Notation for this problem does not usually include the explicit dependence on the assignment algorithm, masking its contribution to the uncertainty in the whole procedure. The choice of Algorithm is left to the experimenter, which means  $z_{\Omega_{\text{Alg}}}$  depends on arbitrarily chosen aspects like the order of the data, and  $\text{P-value}_{\Omega_{\text{Alg}}}$  suffers the same fate. Partly because of this, one cannot truly study its (finite sample) properties. Rosenbaum (2012) warns that one cannot look at many sets of pairs and pick the one providing the conclusion that the experimenter desires. However, the experimenter's assignment algorithm may be biased for their desired conclusion (or away from that conclusion) without the experimenter knowing it. The uncertainty in the assignment procedure is clearly *not* necessarily a random source of bias.

Moreover:

1. In the classical paradigm, the experimenter does not usually consider what bias is caused by a particular optimization method for the assignment.
2. We have no way of quantifying the uncertainty that comes from the matching procedure if the algorithm is chosen arbitrarily by the experimenter. We must quantify this in order for the result to be robust to the choice of experimenter.
3. It is not our desire to place a probability distribution over the ways that human experimenters choose assignment algorithms. This is not interesting. Instead, what is interesting is the range of results that reasonable experimenters might obtain.

4. What we will propose is not equivalent to taking random subsamples of data and repeating the experiment. That procedure could yield trivial (and disastrous) results in the case where one calculates the maximum and minimum over test statistics. This range would grow infinitely large as the number of observations increases. In contrast, our range would generally decrease as the number of observations grows. We assume all observations are available to be matched, and we robustify only over the assignment procedure.

## 2.2. Proposed Approach

First we define a *set of good assignments* as  $\mathcal{A}_{\text{good}} \subseteq \mathcal{A}$  where all  $\Omega \in \mathcal{A}_{\text{good}}$  obey constraints besides those in  $\mathcal{A}$  such as (for example):

- (Calipers) When  $\Omega(i) \neq \emptyset$  then  $\text{dist}(\mathbf{x}_i^t, \mathbf{x}_{\Omega(i)}^c) \leq \epsilon$ .
- (Covariate balance, mean of chosen treatment units close to mean of full treatment group)

$\forall$  covariates  $p$ , using notation  $\bar{x}_p^t = \frac{1}{T} \sum_{t=1}^T x_{ip}^t$ ,

$$\frac{1}{|\Omega|} \sum_{\{i:\Omega(i) \neq \emptyset\}} x_{ip}^t - \bar{x}_p^t \leq \epsilon_p \text{ and}$$

$$\bar{x}_p^t - \frac{1}{|\Omega|} \sum_{\{i:\Omega(i) \neq \emptyset\}} x_{ip}^t \leq \epsilon_p.$$

- (Covariate balance, mean of chosen treatment units similar to mean of full sample)

$\forall$  covariates  $p$ , using notation:

$$\bar{x}_p^{tc} := \frac{1}{T+C} \left( \sum_{t=1}^T x_{ip}^t + \sum_{c=1}^C x_{ip}^c \right), \quad (1)$$

$$\frac{1}{|\Omega|} \sum_{\{i:\Omega(i) \neq \emptyset\}} x_{ip}^t - \bar{x}_p^{tc} \leq \epsilon_p, \text{ and vice versa}$$

$$\bar{x}_p^{tc} - \frac{1}{|\Omega|} \sum_{\{i:\Omega(i) \neq \emptyset\}} x_{ip}^t \leq \epsilon_p.$$

- (Maximizing the fitness of the matches) One can optimize a measure of user-defined fitness for the assignment, and then constrain  $\mathcal{A}_{\text{good}}$  to include all other feasible assignments at or near that fitness level, by including the following constraints:

$$\text{Fitness}(\Omega, \{x_i^t\}_{i=1}^T, \{x_i^c\}_{i=1}^C) \geq \text{Maxfit} - \epsilon,$$



where Maxfit is precomputed,

$$\text{Maxfit} = \max_{\Omega' \in \mathcal{A}} \text{Fitness}(\Omega', \{x_i^t\}_{i=1}^T, \{x_i^c\}_{i=1}^C).$$

If one desires the range of results for all maximally fit pairs and no other pairs,  $\epsilon$  can be set to 0. In our examples we use calipers mostly for ease of notation, but replacing constraints with those above (or other constraints) is trivially simple using MIP software.

$\mathcal{A}_{\text{good}}$  intuitively represents the set of assignments arising from all reasonable matching methods.

The procedure is as follows:

**Robust Procedure.** Define  $\mathcal{A}_{\text{good}}$ . Compute and return the maximum and minimum P-value over the set  $\mathcal{A}_{\text{good}}$ , the maximum and minimum value of the test statistic, and the maximum and minimum value of the estimated treatment effect.

For the matched pair  $z$ -test we would write:

$$\text{P-value}_{\max} = \max_{\Omega \in \mathcal{A}_{\text{good}}} 1 - \Phi(z_{\Omega}) = 1 - \Phi\left(\min_{\Omega \in \mathcal{A}_{\text{good}}} z_{\Omega}\right)$$

and

$$\text{P-value}_{\min} = \min_{\Omega \in \mathcal{A}_{\text{good}}} 1 - \Phi(z_{\Omega}) = 1 - \Phi\left(\max_{\Omega \in \mathcal{A}_{\text{good}}} z_{\Omega}\right).$$

The two P-values quantify the possible uncertainty due to the matching procedure. If all reasonable matching algorithms for the Classical Procedure produce an  $\Omega_{\text{Alg}} \in \mathcal{A}_{\text{good}}$ , this would imply:

$$\text{For all } \Omega_{\text{Alg}} \in \mathcal{A}_{\text{good}}, \text{P-value}_{\min} \leq \text{P-value}_{\Omega_{\text{Alg}}} \leq \text{P-value}_{\max}.$$

Thus, computing  $\text{P-value}_{\min}$  and  $\text{P-value}_{\max}$  would be robust to the choices of human behavior that influence the algorithm for choosing  $\Omega_{\text{Alg}}$ .

The procedure also returns the range of  $z$ -scores,  $\min_{\Omega \in \mathcal{A}_{\text{good}}} z_{\Omega}$  and  $\max_{\Omega \in \mathcal{A}_{\text{good}}} z_{\Omega}$ .

The procedure returns the range of causal effects, which are

$$\left[ \max_{\Omega \in \mathcal{A}_{\text{good}}, \sigma_{\Omega}^2 < b} \bar{d}_{\Omega}, \min_{\Omega \in \mathcal{A}_{\text{good}}, \sigma_{\Omega}^2 < b} \bar{d}_{\Omega} \right], \text{ where } \bar{d}_{\Omega} = \frac{1}{|\Omega|} \sum_{\{i: \Omega(i) \neq \emptyset\}} (y_i^t - y_{\Omega(i)}^c).$$

Here the constraint  $\sigma_{\Omega}^2 < b$  ensures that the causal effect has a sufficiently small variance among the chosen matches, where  $\sigma_{\Omega}^2 = \text{Var}(\{y_i - y_{\Omega(i)}\}_{\{i: \Omega(i) \neq \emptyset\}})$ . This ensures the range is meaningful.

Extending our reasoning from bullet 3 in Section 2.1,

- We do not want to consider how likely a certain assignment  $\Omega$  is to appear. This would involve modeling human behavior of the experimenter. We do not want to place a distribution over the choice over algorithms that an experimenter would choose. As Morgan and Winship (2007) note, there is no clear guidance on the choice of matching procedure. We do not presuppose a distribution over these procedures.

- We do not want to consider statistics of the set of  $\mathcal{A}_{\text{good}}$ , such as the average P-value over the set of  $\mathcal{A}_{\text{good}}$ . This is simply a special case of the previous point, where we assume that all good assignments are equally likely to be chosen, which is clearly not the case.

A remark on the covariate balance constraints: as noted in many other works, by choosing constraints on  $\Omega$  and/or an objective that favors using more treatment points, the estimates of causal effect will generally be less biased for the ATT. One can similarly alter constraints or the objective to reduce bias on ATE. Many other works speak to the benefits of better matching procedures and reduction of bias (e.g., Rosenbaum 1989).

A remark on the uses of this procedure: the techniques proposed here are appropriate for observational data, and not appropriate for studies where the matching is done prior to data collection, for the reason that one generally would not want to match, collect data, then afterwards rethink the assignments among the collected data.

In what follows, we provide special cases of the Robust Procedure for various specific hypothesis tests. The goal is to provide conclusions that are robust to the class of experimenters within  $\mathcal{A}_{\text{good}}$ .

### 3. A Robust $Z$ -Test

Let us consider the upper one-sided  $z$ -test for estimating whether the difference in mean of the treatment and control populations is sufficiently greater than 0. Formulations that are slightly different from each other (that we will discuss) can handle the ATT estimation problem, the ATE problem, or the overlap problem of Rosenbaum (2012). The only difference between these formulations is in the constraints.

This test will also encompass the upper one-sided  $t$ -test when the sample size is sufficiently large ( $\geq 30$ ). (Because of the fast convergence of the sample standard deviation to the true standard

deviation, in practice we use the sample standard deviation for the calculations that follow, as is usual practice.) From each pair of observations, we record  $y_i^t - y_{\Omega(i)}^c$ , which is the difference between the treatment and control outcomes for pair  $i$ , that is, (treatment outcome)-(control outcome). The average estimated treatment effect on the samples is  $\bar{d}_\Omega$ , which is the sample average of the differences,  $\bar{d}_\Omega = \frac{1}{|\Omega|} \sum_{\{i:\Omega(i) \neq \emptyset\}} y_i^t - y_{\Omega(i)}^c$ . The mean of the treatment and control populations for the high density shared population region are denoted by  $\mu_T$  and  $\mu_C$ , and the treatment effect is  $\mu_D = \mu_T - \mu_C$ . We would like to know whether  $\mu_D$  is greater than 0, meaning that the treatment has an effect. We consider the hypothesis  $H_0 : \mu_D = 0$  versus the alternative  $H_1 : \mu_D > 0$ . In what follows,  $\Phi$  is the cumulative density function of a standard normal distribution,  $n$  is the total number of pairs, and  $\hat{\sigma}$  is the sample standard deviation of the differences,  $y_i^t - y_{\Omega(i)}^c$ . The  $z$ -score is:

$$z = \frac{\bar{d}_\Omega \sqrt{n}}{\hat{\sigma}},$$

and the  $P$ -value is

$$P\text{-value} = 1 - \Phi(z).$$

Since  $1 - \Phi(z)$  is monotonically decreasing in  $z$ , maximizing  $z$  is the same as minimizing the value of  $1 - \Phi(z)$ .

We now provide an algorithm to choose matched pairs in order to maximize  $z$  (and we obtain  $z_{\max}$  where  $z_{\max} = \max_{\mathbf{a} \in \mathcal{A}_{\text{good}}} z(\mathbf{a})$ ), and a simple adaptation can be made to minimize  $z$  (and we obtain  $z_{\min}$  where  $z_{\min} = \min_{\mathbf{a} \in \mathcal{A}_{\text{good}}} z(\mathbf{a})$ ). As discussed above, the quantities we aim to report are the ranges of  $z$ -scores and  $P$ -values, that is,  $[z_{\min}, z_{\max}]$ , and  $[1 - \Phi(z_{\max}), 1 - \Phi(z_{\min})]$ . Along the way, we will derive a formulation to report the range of causal effects.

### 3.1. Formulations

The input parameters are as follows:

$n$  is the total numbers of matched pairs we will find. For the full matching ATT problem, we set  $n = T$ , where  $T$  is the total number of treatment points in the data.

$Q$  is the set of all observations in the treatment group, indexed by  $i$

$R$  is the set of all observations in the control group, indexed by  $j$

$T_i$  is the outcome of a treated observation  $i$  in the treatment group

$C_j$  is the outcome of a control observation  $j$  in the control group

$\text{dist}_{ij}$  is the  $ij$ th element of a matrix. It takes value 1 if the covariates of treated observation  $i$  and control observation  $j$  are similar enough to be a possible matched pair, otherwise 0.

$\hat{\sigma}$  is the sample standard deviation of the paired differences, which closely approximates the true standard deviation when the number of observations is sufficiently large. We use the uncorrected standard deviation here for ease of exposition.

There is only one matrix of decision variables, namely:

$a_{ij}$  is a binary variable that is 1 if  $i$  and  $j$  are in the same pair, otherwise 0.

We start with the ATT problem with full matching on the treatment cases. Optimization for  $z = \bar{d}/(\hat{\sigma}/\sqrt{n})$  can be formulated as follows.

**Formulation 1 (Nonlinear  $z$ -test for ATT with full matching)**

$$\text{Maximize/Minimize}_{\mathbf{a}} \quad z(\mathbf{a}) = \frac{\frac{1}{T} \sum_{i \in Q} \sum_{j \in R} (T_i - C_j) a_{ij} \sqrt{T}}{\hat{\sigma}}$$

subject to:

$$\begin{aligned} \hat{\sigma} &= \sqrt{\frac{1}{T} \sum_{i \in Q} \sum_{j \in R} [(T_i - C_j) a_{ij}]^2 - \left( \frac{1}{T} \sum_{i \in Q} \sum_{j \in R} (T_i - C_j) a_{ij} \right)^2} \\ \sum_{i \in Q} \sum_{j \in R} a_{ij} &= T && \text{(Choose } T \text{ pairs, all treatment points are matched)} \\ \sum_{i \in Q} a_{ij} &\leq 1 && \forall j \quad \text{(Choose at most one treatment observation)} \\ \sum_{j \in R} a_{ij} &\leq 1 && \forall i \quad \text{(Choose at most one control observation)} \\ a_{ij} &\leq \text{dist}_{ij} && \forall i, j \quad \text{(Choose only pairs that are allowed)} \\ a_{ij} &\in \{0, 1\} && \forall i, j \quad \text{(Defines binary variable } a_{ij}) \\ &&& \text{(Additional user-defined covariate balance constraints.)} \end{aligned}$$

The term in the objective function  $\frac{1}{n} \sum_{i \in Q} \sum_{j \in R} (T_i - C_j) a_{ij}$  is the value of  $\bar{d}$ , where the  $a_{ij}$  term ensures that only differences between the pairs  $i, j$  that we are selecting are used for calculating  $\bar{d}$ . The first constraint ensures that we choose exactly all  $T$  pairs, the second and third constraints ensure that only one treatment and one control observation are selected for each pair, and the last constraint is an if-then constraint, stating that we are only allowed to choose pairs  $i, j$  for which  $\text{dist}_{ij} = 1$ . Recall that the  $\text{dist}_{ij}$ 's were determined in advance, and they encode whether  $i$ 's covariates are close enough to  $j$ 's covariates to be chosen as a possible pair.

Formulation 1 can be modified to form the problem of choosing both treatment and control populations simultaneously, to handle the setting of Rosenbaum (2012), which we do next. Here, the mean is taken over the same portion of the population as Rosenbaum (2012), which is the region of overlap between the control and treatment populations, removing extreme regions. The number of pairs will be fixed at  $n$  and we loop over all values of  $n$  where feasible solutions exist.

**Formulation 2 (Nonlinear  $z$ -test for Overlap Problem)**

$$\text{Maximize/Minimize}_{\mathbf{a}} \quad z(\mathbf{a}) = \frac{\frac{1}{n} \sum_{i \in Q} \sum_{j \in R} (T_i - C_j) a_{ij} \sqrt{n}}{\hat{\sigma}}$$

subject to:

$$\begin{aligned} \hat{\sigma} &= \sqrt{\frac{1}{n} \sum_{i \in Q} \sum_{j \in R} [(T_i - C_j) a_{ij}]^2 - \left( \frac{1}{n} \sum_{i \in Q} \sum_{j \in R} (T_i - C_j) a_{ij} \right)^2} \\ \sum_{i \in Q} \sum_{j \in R} a_{ij} &= n && \text{(Choose } n \text{ pairs)} \\ \sum_{i \in Q} a_{ij} &\leq 1 && \forall j \quad \text{(Choose at most one treatment observation)} \\ \sum_{j \in R} a_{ij} &\leq 1 && \forall i \quad \text{(Choose at most one control observation)} \\ a_{ij} &\leq \text{dist}_{ij} && \forall i, j \quad \text{(Choose only pairs that are allowed)} \\ a_{ij} &\in \{0, 1\} && \forall i, j \quad \text{(Defines binary variable } a_{ij}) \\ &&& \text{(Additional user-defined covariate balance constraints.)} \end{aligned}$$

**Covariate Balance Constraints.** For any of the formulations we provide in this work, one can replace or augment the  $a_{ij} \leq \text{dist}_{ij}$  constraints with more sophisticated distributional matching

constraints such as those of Zubizarreta (2012), Zubizarreta et al. (2013) or Zubizarreta et al. (2014), including the constraints suggested earlier. To turn this into a  $z$ -test for the ATT problem, one can either add balance constraints that force the chosen treatment and control patients to have similar statistics to the full treatment group, or simplify the problem substantially by forcing all treatment points to be used. For example, we might add the following covariate balance constraints to the formulation, which could be used for the ATE problem, encoding (1). For each covariate  $p$ , we would include the following two constraints into our formulation, which are linear in the decision variables  $a_{ij}$ .

$$\begin{aligned}\frac{1}{n} \sum_{j \in R} \sum_{i \in Q} x_{ip}^t a_{ij} - \bar{x}_p^{tc} &\leq \epsilon_p, \\ \bar{x}_p^{tc} - \frac{1}{n} \sum_{i \in Q} \sum_{j \in R} x_{ip}^t a_{ij} &\leq \epsilon_p,\end{aligned}$$

and similarly for the control group:

$$\begin{aligned}\frac{1}{n} \sum_{j \in R} \sum_{i \in Q} x_{jp}^c a_{ij} - \bar{x}_p^{tc} &\leq \epsilon_p, \\ \bar{x}_p^{tc} - \frac{1}{n} \sum_{i \in Q} \sum_{j \in R} x_{jp}^c a_{ij} &\leq \epsilon_p.\end{aligned}$$

### 3.2. Linear Formulations $z$ -test

The formulations above are clearly not linear in the decision variables. Let us consider the solution of Formulation 2, since Formulation 1 is a special case. Its solution can be approximated using a MINLP (mixed-integer nonlinear programming solver) but no guarantees on the optimality of the solution can be made. In what follows, we show how this problem can be simplified to be solved by an algorithm that solves several linear integer programming problems instead. This algorithm benefits from the computational speed of ILP solvers and has a guarantee on the optimality of the solution.

To create the ILP formulation, we note that the objective is increasing in the average of the differences (this term appears both in the numerator and denominator), and it is decreasing in

the sum of the squared differences (this term is the first term of  $\hat{\sigma}$ ). We will replace the nonlinear objective of Formulation 2 as follows:

$$\text{Maximize/Minimize}_{\mathbf{a}} \sum_{i \in Q} \sum_{j \in R} (T_i - C_j) a_{ij}, \quad (2)$$

which is now linear. The quantity in (2) is the estimated causal effect. Thus we will be finding the range of causal effects.

At the same time, we will limit the sum of squared differences term by  $b_l$ , which is now a parameter rather than a decision variable. Thus, we will optimize causal effect subject to a bound on the variance. We introduce the following new constraint:

$$\sum_{i \in Q} \sum_{j \in R} [(T_i - C_j) a_{ij}]^2 \leq b_l. \quad (3)$$

We will show in Section 3.2.1 how to choose  $b_l$ 's in order to maintain the guarantee of optimality of the solution. Since  $a_{ij}$  is a binary variable, the nonlinear constraint (3) can be replaced by the following equivalent linear constraint:

$$\sum_{i \in Q} \sum_{j \in R} (T_i - C_j)^2 a_{ij} \leq b_l. \quad (4)$$

Putting this together, the new formulation is an ILP.

**Formulation 3 (Optimize causal effect for  $z$ -test, with upper bound on variance)**

$$\text{Maximize/Minimize}_{\mathbf{a}} \sum_{i \in Q} \sum_{j \in R} (T_i - C_j) a_{ij} \quad (\text{Causal effect})$$

subject to:

$$\begin{aligned} \sum_{i \in Q} \sum_{j \in R} (T_i - C_j)^2 a_{ij} &\leq b_l && (\text{Upper bound on sample variance}) \\ \sum_{i \in Q} \sum_{j \in R} a_{ij} &= n && (\text{Choose } n \text{ pairs}) \\ \sum_{i \in Q} a_{ij} &\leq 1 && \forall j \quad (\text{Choose at most one treatment observation}) \\ \sum_{j \in R} a_{ij} &\leq 1 && \forall i \quad (\text{Choose at most one control observation}) \\ a_{ij} &\leq \text{dist}_{ij} && \forall i, j \quad (\text{Choose only pairs that are allowed}) \\ a_{ij} &\in \{0, 1\} && \forall i, j \quad (\text{Defines binary variable } a_{ij}) \\ &&& (\text{Additional user-defined covariate balance constraints.}) \end{aligned}$$

This formulation optimizes causal effect, subject to the variance of the causal effect being small. This formulation can be used by itself to find the range of reasonable causal effects, given a fixed bound  $b_l$  on the variance.

Let us get back to optimizing the z-score. Our algorithm will solve this formulation for many different values of  $b_l$  to find the optimal z-score and P-value. If we denote the solution of the maximization problem as  $a_l$ , where  $a_l$  is still also indexed by  $ij$ , we will then be able to bound the value of  $z$  (as we will show through Theorem 1 later). To write the notation for the bound, we define:

$$\bar{d}_{a_l} := \frac{1}{n} \sum_{i \in Q} \sum_{j \in R} (T_i - C_j) a_{l,ij}. \quad (5)$$

Shortly we will use Theorem 1 below to prove an upper bound on the z-score as follows:

$$\max_{\mathbf{a}} z(\mathbf{a}) \leq \frac{\bar{d}_{a_{l^*}} \sqrt{n}}{\sqrt{\frac{1}{n} b_{l^*-1} - (\bar{d}_{a_{l^*}})^2}} \leq \max_l \frac{\bar{d}_{a_l} \sqrt{n}}{\sqrt{\frac{1}{n} b_{l-1} - (\bar{d}_{a_l})^2}}, \quad (6)$$

where  $l^*$  is the index where

$$b_{l^*-1} \leq \sum_{i \in Q} \sum_{j \in R} (T_i - C_j)^2 a_{ij}^* \leq b_{l^*}.$$

There will be an analogous lower bound from Theorem 1, which is:

$$\max_{\mathbf{a}} z(\mathbf{a}) \geq \max_l \frac{\bar{d}_{a_l} \sqrt{n}}{\sqrt{\frac{1}{n} b_l - (\bar{d}_{a_l})^2}}, \quad (7)$$

The  $P$ -value is  $1 - \Phi(z)$ , so upper and lower bounds for  $z$  directly yield lower and upper bounds for the  $P$ -value.

In what follows we show how to choose progressively finer meshes for  $b_l$  to maintain the guarantee on the quality of the solution, repeatedly solving Formulation 3.

**3.2.1. Algorithm for  $z$ -test with performance guarantees** Let us first derive the upper and lower bounds in (6) and (7). They come directly from the following:

**THEOREM 1.** *Consider the optimization problem*

$$x^* \in \arg \max_x F(f_1(x), f_2(x)),$$



where  $f_1$  and  $f_2$  are real-valued functions of  $x \in X$ ,  $F$  is monotonically increasing in  $f_1$  and monotonically decreasing in  $f_2$ . Assume we are given  $[b_1, b_2, \dots, b_l, \dots, b_L]$  that span a wide enough range so that  $x^*$  obeys:

$$b_{l^*-1} \leq f_2(x^*) \leq b_{l^*} \text{ for some } l^* \in \{1, \dots, L\}.$$

Define  $x_l$  as follows:

$$x_l \in \arg \max_{x: f_2(x) \leq b_l} F(f_1(x), b_l) = \arg \max_{x: f_2(x) \leq b_l} f_1(x),$$

where the equality follows because  $F$  monotonically increases in  $f_1$ . Then

$$\max_l F(f_1(x_l), b_l) \leq \max_x F(f_1(x), f_2(x)) \leq F(f_1(x_{l^*}), b_{l^*-1}) \leq \max_l F(f_1(x_l), b_{l-1}).$$

This theorem bounds the optimal value of  $F$  along the whole regime of  $x$  in terms of the values computed at the  $L$  grid points. The proof of Theorem 1 is not difficult and can be found in Section 3.3.

Note that the objective function of Formulations 1 and 2 is exactly of the form of Theorem 1, where

$$f_1(\mathbf{a}) = \frac{1}{n} \sum_{i \in Q} \sum_{j \in R} (T_i - C_j) a_{ij},$$

$$f_2(\mathbf{a}) = \sum_{i \in Q} \sum_{j \in R} (T_i - C_j)^2 a_{ij},$$

and

$$F(f_1(\mathbf{a}), f_2(\mathbf{a})) = \frac{f_1(\mathbf{a})\sqrt{n}}{\sqrt{\frac{1}{n}f_2(\mathbf{a}) - (f_1(\mathbf{a}))^2}}.$$

The extra constraints on  $\mathbf{a}$  in Formulation 1 are also compatible with Theorem 1. Thus, the bounds (6) and (7) are direct results of Theorem 1 applied to Formulation 1's objective function.

Now that this is established, we write out the algorithm, starting by maximizing the  $z$ -score. The minimization algorithm is symmetric to the maximization algorithm.

#### Algorithm 1: Maximize $z$ -score

*Step 1: (Compute the value for  $b_L$ )* We solve Formulation 3 by relaxing (removing) the

first constraint (upper bound on sample standard deviation) and compute the value of  $b_L$  as

$$\sum_{i \in Q} \sum_{j \in R} (T_i - C_j)^2 a_{ij}.$$

*Step 2: (Create Coarse Mesh)* We create a coarse mesh  $b_1, \dots, b_L$  where  $b_1 < b_l < b_L$ . We want the interval  $[b_1, b_L]$  to be wide enough to contain the true value of  $f_2(\mathbf{a}^*) := \sum_{i \in Q} \sum_{j \in R} (T_i - C_j)^2 a_{ij}^*$ , where  $\mathbf{a}^* \in \arg \max z(\mathbf{a})$ , which we do not know and are trying to estimate. Usually we choose the  $b_l$  evenly spaced, though they do not need to be.

*Step 3: (Solve ILP's Along Mesh)* For each  $l$ , we compute the solution to the ILP Formulation 2.

We then compute upper and lower bounds for the solution using (6) and (7). In particular:

$$\max_l \frac{\bar{d}_{a_l} \sqrt{n}}{\sqrt{\frac{1}{n} b_l - (\bar{d}_{a_l})^2}} \leq \max_a z(\mathbf{a}) \leq \max_l \frac{\bar{d}_{a_l} \sqrt{n}}{\sqrt{\frac{1}{n} b_{l-1} - (\bar{d}_{a_l})^2}}. \quad (8)$$

*Step 4: (Exclude Irrelevant Mesh Intervals)* For each  $l$  we determine whether the interval  $[b_{l-1}, b_l]$  can be excluded because it provably does not contain a  $f_2(\mathbf{a})$  value corresponding to the maximum value of  $z(\mathbf{a})$ . In particular, we know from the bounds (7) and (6) that if the upper bound on the objective for a particular  $b_{l'}$  is lower than all lower bounds for the optimal solution  $\mathbf{a}^*$  then  $l'$  cannot equal  $l^*$  and the interval  $[b_{l'-1}, b_{l'}]$  can be excluded from further exploration. Specifically, we check for each  $l$  whether

$$\frac{\bar{d}_{a_l} \sqrt{n}}{\sqrt{\frac{1}{n} b_{l-1} - (\bar{d}_{a_l})^2}} < \max_l \frac{\bar{d}_{a_l} \sqrt{n}}{\sqrt{\frac{1}{n} b_l - (\bar{d}_{a_l})^2}}.$$

If this holds for some  $l$ , it means  $l$  cannot equal  $l^*$  and the interval  $[b_{l-1}, b_l]$  can be excluded from further exploration. We create a set of the intervals  $[b_{l-1}, b_l]$  that were not excluded by this process.

*Step 5: (Refine Mesh)* Create a finer mesh by redefining the  $b_l$ 's to be closer together within the intervals that are still included. This will yield tighter upper and lower bounds. We then re-solve the ILP's for every new  $b_l$  and compute the bounds (8).

*Step 6: (Repeat)* Repeat Step 4, which is to exclude irrelevant intervals, and Step 5, which is to refine the mesh and re-solve the ILP's, until desired tolerance  $\epsilon$  is reached in the bounds (8):

$$\max_l \frac{\bar{d}_{a_l} \sqrt{n}}{\sqrt{\frac{1}{n} b_{l-1} - (\bar{d}_{a_l})^2}} - \max_l \frac{\bar{d}_{a_l} \sqrt{n}}{\sqrt{\frac{1}{n} b_l - (\bar{d}_{a_l})^2}} \leq \epsilon.$$

### 3.3. Proof of Theorem 1

Because the  $b_l$ 's are defined on a pre-specified grid, we know the maximum value of  $F$  may not occur at one of the grid points. Since by definition of  $x_l$  we know that  $f_2(x_l) \leq b_l$ , and since  $F$  is decreasing in its second argument, we have  $F(f_1(x_l), b_l) \leq F(f_1(x_l), f_2(x_l))$  for each  $l$ , and taking a max over all  $l$ :

$$\max_{l \in 1 \dots L} F(f_1(x_l), b_l) \leq \max_{l \in 1 \dots L} F(f_1(x_l), f_2(x_l)) \leq \max_x F(f_1(x), f_2(x)).$$

This is the left inequality of the bound. The rest of the proof deals with the right inequalities.

First, it is true that:

$$f_1(x^*) = \max_{x: f_2(x) = f_2(x^*)} f_1(x). \quad (9)$$

If this were not true then either  $f_1(x^*) < \max_{x: f_2(x) = f_2(x^*)} f_1(x)$  or  $f_1(x^*) > \max_{x: f_2(x) = f_2(x^*)} f_1(x)$ . If the first inequality were true then

$$F(f_1(x^*), f_2(x^*)) < F\left(\max_{x: f_2(x) = f_2(x^*)} f_1(x), f_2(x^*)\right),$$

which contradicts the definition of  $x^*$ . The second option also cannot be true as we know there exists a solution  $x^*$  so that the maximum is attained with  $f_1$  and  $f_2$  values  $f_1(x^*)$  and  $f_2(x^*)$ . So we can say that (9) holds.

From (9) and using  $l^*$  defined in the statement of the theorem, we can derive:

$$f_1(x^*) = \max_{x: f_2(x) = f_2(x^*)} f_1(x) \leq \max_{x: f_2(x) \leq f_2(x^*)} f_1(x) \leq \max_{x: f_2(x) \leq b_{l^*}} f_1(x) = f_1(x_{l^*}), \quad (10)$$

where we used that the set  $\{x: f_2(x) = f_2(x^*)\}$  is smaller than the set  $\{x: f_2(x) \leq f_2(x^*)\}$  which is smaller than  $\{x: f_2(x) \leq b_{l^*}\}$ , since  $f_2(x^*) \leq b_{l^*}$  by definition of  $l^*$ . Thus,  $f_1(x^*) \leq f_1(x_{l^*})$ . Now,

$$\begin{aligned} F(f_1(x^*), f_2(x^*)) &\leq F(f_1(x^*), b_{l^*-1}) \\ &\leq F(f_1(x_{l^*}), b_{l^*-1}) \\ &\leq \max_l F(f_1(x_l), b_{l-1}). \end{aligned}$$

Here the first inequality above follows from the definition of  $l^*$ ,  $b_{l^*-1} \leq f_2(x^*)$ , and the fact that  $F$  decreases in the second argument. The second inequality comes from (10) and the fact that  $F$  is increasing in its first argument. The third inequality follows from taking a maximum over all  $l$  rather than using  $l^*$ . The proof is complete.

The results presented so far are only suitable for testing data associated with real-valued outcomes. In order to test count/proportion data (i.e., with binary outcomes), we need to consider a robust version of McNemar’s test.

#### 4. Robust McNemar’s Test

Let us consider casual inference for binary outcomes. Here, we will use McNemar’s test for matched pairs with binary outcomes. Again we will consider the overlap problem, where the distribution being considered is the overlap between high density regions of the treatment and control populations. Again for the ATT problem we could introduce balance constraints, or simplify the problem by forcing the formulation to use all treatment units. The  $P$ -value is  $1 - \Phi(z)$  for sufficiently large samples, and  $z$  can be calculated by the following equation, where  $B$  and  $C$  are counts of discordants/untied responses:

$$z = \left[ \frac{B - C - 1}{\sqrt{B + C}} \right]. \quad (11)$$

Here, one can think of  $B$  as the number of pairs where the outcome from the treated patient was “Yes” and the outcome for the untreated patient was “No”.  $C$  is the number of pairs where the outcome from the treated patient was “No” and the outcome from the untreated patient was “Yes”.

The formulation below optimizes the causal effect, measured by how much larger the untied responses  $B$  are than the untied responses  $C$ , relative to the total number of untied responses  $B + C$ . The formulation also maximize/minimizes the  $P$ -value, by minimizing/maximizing  $z$ . We define the following parameters to formulate the model:

$m$  is the total numbers of untied responses. We loop over all possible values of  $m$ , until the solution becomes infeasible.

$T_i$  is the outcome of a treated observation  $i$  in the treatment group

$C_j$  is the outcome of a control observation  $j$  in the control group

$\text{dist}_{ij}$  is the  $ij$ th element of a matrix. It takes value 1 if the covariates of treated observation  $i$  and control observation  $j$  are similar enough to be a possible matched pair, otherwise 0.

$b_{ij}$  is 1 if  $a_{ij}$  and  $C_j$  are equal to 1 and  $T_i$  is equal to 0, otherwise 0 (first type of discordant pair)

$c_{ij}$  is 1 if  $a_{ij}$  and  $T_i$  are equal to 1 and  $C_j$  is equal to 0, otherwise 0 (second type of discordant pair)

There is only one matrix of decision variables, namely:

$a_{ij}$  is a binary variable that is 1 if  $i$  and  $j$  are in the same pair, otherwise 0.

One can show that the number of pairs where the same outcome is realized for treatment and control is irrelevant, so we allow it to be chosen arbitrarily, with no constraints or variables defining it. The total number of pairs  $n$  is also not relevant for this test. Therefore, we choose only the total number of untied responses  $m$ . Finally, the formulation becomes:

$$\text{Maximize/Minimize}_{\mathbf{a}} \quad z(\mathbf{a}) = \left\lceil \frac{B - C - 1}{\sqrt{B + C}} \right\rceil$$

subject to:

$$b_{ij} = a_{ij}C_j(1 - T_i) \quad \forall i, j \quad (\text{Defines } b_{ij}) \quad (12)$$

$$c_{ij} = a_{ij}T_i(1 - C_j) \quad \forall i, j \quad (\text{Defines } c_{ij}) \quad (13)$$

$$\sum_{i \in Q} \sum_{j \in R} b_{ij} = B \quad (\text{Total number of first type of discordant pairs}) \quad (14)$$

$$\sum_{i \in Q} \sum_{j \in R} c_{ij} = C \quad (\text{Total number of second type of discordant pairs}) \quad (15)$$

$$B + C = m \quad (\text{Total number of discordant pairs}) \quad (16)$$

$$\sum_{i \in Q} a_{ij} \leq 1 \quad \forall j \quad (\text{Choose at most one treatment observation}) \quad (17)$$

$$\sum_{j \in R} a_{ij} \leq 1 \quad \forall i \quad (\text{Choose at most one control observation}) \quad (18)$$

$$a_{ij} \leq \text{dist}_{ij} \quad \forall i, j \quad (\text{Choose only pairs that are allowed}) \quad (19)$$

$$a_{ij} \in \{0, 1\} \quad \forall i, j \quad (\text{Defines binary variable } a_{ij}) \quad (20)$$

$$(\text{Additional user-defined covariate balance constraints.}) \quad (21)$$

Equations (12) and (13) define variables  $b_{ij}$  and  $c_{ij}$ . Equations (14) and (15) are used to define variables  $B$  and  $C$ . To control the total number of untied responses, we incorporate Equation (16). Equations (17) and (18) confirm that only one treated/control unit will be assigned in a single pair. Equation (19) says that the variable  $a_{ij}$  can take value 1 only if the value of parameter  $\text{dist}_{ij}$  is 1. If we add the constraint  $\sum_{i \in Q} \sum_{j \in R} a_{ij} = T$ , where  $T$  is the total number of treatment points, we will be in the case of estimating the ATT using all treatment points.

This test is equivalent to a particular  $\chi^2$  test (see Tamhane and Dunlop 2000).

## 5. A Robust $\chi^2$ Test

We generalize the previous result to consider 1:M matching, i.e., more than one control units are going to be matched with 1 treatment unit. Here, we will use the  $\chi^2$  test for 1:M matched pairs with binary outcomes. For introductory material, see Breslow and Day (1980). The relevant test statistic is as follows, where  $T_{m,M} = n_{1,m-1,M} + n_{0,m,M}$ ,  $n_{1,m-1,M}$  denotes the number of matched sets containing  $M$  control units and one treatment unit, where the treatment unit has positive outcome and  $m-1$  of the control units have positive outcome. Let  $n_{0,m,M}$  denotes the number of matched sets containing  $M$  control units of which  $m$  have positive outcome and where the outcome of the treated unit is 0.

$$\chi^2 = \frac{\left[ \sum_M \sum_{m=1}^M \left( n_{1,m-1,M} - \frac{m}{M+1} T_{m,M} \right) - 1/2 \right]^2}{\sum_M \sum_{m=1}^M T_{m,M} \frac{m(M-m+1)}{(M+1)^2}}. \quad (22)$$

In order to optimize  $\chi^2$  using an ILP technique, we will optimize the value of  $\chi$  and then calculate the  $\chi^2$  value. We will use the following variables :

According to the definition of  $n_{1,m-1,M}$ ,  $n_{1,m-1,M} = \sum_i T_i \Omega_{i,m-1,M} \beta_{i,m-1,M}$ , where  $\Omega_{i,m-1,M}$  and  $\beta_{i,m-1,M}$  are binary variables with the following definitions: If  $\sum_j a_{ij} = M$  then  $\Omega_{i,m-1,M} = 1$ , 0 otherwise; If  $\sum_j a_{ij} C_j = m-1$  then  $\beta_{i,m-1,M} = 1$ , 0 otherwise.

Similarly,  $n_{0,m,M} = \sum_i (1 - T_i) \theta_{i,m,M} \gamma_{i,m,M}$ , where  $\theta_{i,m,M}$  and  $\gamma_{i,m,M}$  are binary variables with the following definitions: If  $\sum_j a_{ij} = M$  then  $\theta_{i,m,M} = 1$ , 0 otherwise; If  $\sum_j a_{ij} C_j = m$  then  $\gamma_{i,m,M} = 1$ , 0 otherwise.

In order to develop an equivalent linear model we can introduce two new binary variables  $\lambda_{i,m-1,M}$  and  $\zeta_{i,m,M}$  and, define  $n_{1,m-1,M}$  as  $\sum_i T_i \lambda_{i,m-1,M}$ , where  $\lambda_{i,m-1,M} = \Omega_{i,m-1,M} \beta_{i,m-1,M}$  and  $n_{0,m,M}$  as  $\sum_i (1 - T_i) \zeta_{i,m,M}$ , where  $\zeta_{i,m,M} = \theta_{i,m,M} \gamma_{i,m,M}$ . We will add linear constraints to restore these relationships, which we will discuss in a moment.

The equivalent linear formulation for maximizing the test statistic can be represented as follows, where again a constant term  $\phi$  is used instead of the denominator, and we add an upper bound on the denominator using  $\phi$ . To get a guaranteed optimal solution, we will solve a series of ILP instances by tuning the value of  $\phi$ . In order to handle the absolute term in the numerator (in case of maximization), we can both maximize and minimize the proposed model and select the maximum absolute solution as a final solution.

$$\text{Maximize/Minimize}_{\mathbf{a}} \quad \sum_M \sum_{m=1}^M \left( n_{1,m-1,M} - \frac{m}{M+1} T_{m,M} \right) - 1/2$$

subject to:

$$\sum_M \sum_{m=1}^M T_{m,M} \frac{m(M-m+1)}{(M+1)^2} \leq \phi \quad (23)$$

$$T_{m,M} = n_{1,m-1,M} + n_{0,m,M} \quad \forall m, M \quad (24)$$

$$n_{1,m-1,M} = \sum_i T_i \lambda_{i,m-1,M} \quad \forall m, M \quad (25)$$

$$n_{0,m,M} = \sum_i (1 - T_i) \zeta_{i,m,M} \quad \forall m, M \quad (26)$$

$$\lambda_{i,m-1,M} \leq \Omega_{i,m-1,M} \quad \forall i, m, M \quad (27)$$

$$\lambda_{i,m-1,M} \leq \beta_{i,m-1,M} \quad \forall i, m, M \quad (28)$$

$$\lambda_{i,m-1,M} \geq \Omega_{i,m-1,M} + \beta_{i,m-1,M} - 1 \quad \forall i, m, M \quad (29)$$

$$\zeta_{i,m,M} \leq \theta_{i,m,M} \quad \forall i, m, M \quad (30)$$

$$\zeta_{i,m,M} \leq \gamma_{i,m,M} \quad \forall i, m, M \quad (31)$$

$$\zeta_{i,m,M} \geq \theta_{i,m,M} + \gamma_{i,m,M} - 1 \quad \forall i, m, M \quad (32)$$

$$\sum_{j \in R} a_{ij} \leq M\pi_{i,m-1,M} + (M+1)\Omega_{i,m-1,M} + (U_a + 2)\rho_{i,m-1,M} - 1 \quad \forall i, m, M \quad (33)$$

$$\sum_{j \in R} a_{ij} \geq \pi_{i,m-1,M} + (M+1)\Omega_{i,m-1,M} + (M+2)\rho_{i,m-1,M} - 1 \quad \forall i, m, M \quad (34)$$

$$\pi_{i,m-1,M} + \Omega_{i,m-1,M} + \rho_{i,m-1,M} = 1 \quad \forall i, m, M \quad (35)$$

$$\sum_{j \in R} a_{ij} C_j \leq (m-1)\mu_{i,m-1,M} + m\beta_{i,m-1,M} + (U_c + 2)\tau_{i,m-1,M} - 1 \quad \forall i, m, M \quad (36)$$

$$\sum_{j \in R} a_{ij} C_j \geq \mu_{i,m-1,M} + m\beta_{i,m-1,M} + (m+1)\tau_{i,m-1,M} - 1 \quad \forall i, m, M \quad (37)$$

$$\mu_{i,m-1,M} + \beta_{i,m-1,M} + \tau_{i,m-1,M} = 1 \quad \forall i, m, M \quad (38)$$

$$\sum_{j \in R} a_{ij} \leq M\delta_{i,m,M} + (M+1)\theta_{i,m,M} + (U_a + 2)\epsilon_{i,m,M} - 1 \quad \forall i, m, M \quad (39)$$

$$\sum_{j \in R} a_{ij} \geq \delta_{i,m,M} + (M+1)\theta_{i,m,M} + (M+2)\epsilon_{i,m,M} - 1 \quad \forall i, m, M \quad (40)$$

$$\delta_{i,m,M} + \theta_{i,m,M} + \epsilon_{i,m,M} = 1 \quad \forall i, m, M \quad (41)$$

$$\sum_{j \in R} a_{ij} C_j \leq m\nu_{i,m,M} + (m+1)\gamma_{i,m,M} + (U_c + 2)\eta_{i,m,M} - 1 \quad \forall i, m, M \quad (42)$$

$$\sum_{j \in R} a_{ij} C_j \geq \nu_{i,m,M} + (m+1)\gamma_{i,m,M} + (m+2)\eta_{i,m,M} - 1 \quad \forall i, m, M \quad (43)$$

$$\nu_{i,m,M} + \gamma_{i,m,M} + \eta_{i,m,M} = 1 \quad \forall i, m, M \quad (44)$$

$$\sum_{i \in Q} a_{ij} \leq 1 \quad \forall j \quad (45)$$

$$\sum_{j \in R} a_{ij} \leq M \quad \forall i \quad (46)$$

$$a_{ij} \leq \text{dist}_{ij} \quad \forall i, j \quad (47)$$

$$a_{ij} \in \{0, 1\} \quad \forall i, j \quad (48)$$

$$(\text{Additional user-defined covariate balance constraints.}) \quad (49)$$

Equation (23) imposes an upper bound on the denominator of our objective function. Equations (24)-(26) define decision variables  $T_{m,M}$ ,  $n_{1,m-1,M}$  and  $n_{0,m,M}$ , where Equations (27)-(29) and (30)-(32) are used to linearize the product of  $\Omega_{i,m-1,M}\beta_{i,m-1,M}$  and  $\theta_{i,m,M}\gamma_{i,m,M}$ . Equations (33)-(35), (36)-(38), (39)-(41) and (42)-(44) are used to linearize if-then constraints to define decision variables  $\Omega_{i,m-1,M}$ ,  $\beta_{i,m-1,M}$ ,  $\theta_{i,m,M}$  and  $\gamma_{i,m,M}$  respectively, where  $\pi_{i,m-1,M}$ ,  $\rho_{i,m-1,M}$ ,  $\mu_{i,m-1,M}$ ,  $\tau_{i,m-1,M}$ ,  $\delta_{i,m,M}$ ,  $\epsilon_{i,m,M}$ ,  $\nu_{i,m,M}$  and  $\eta_{i,m,M}$  are binary variables, and  $U_a$  and  $U_c$  are upper bounds on  $\sum_{j \in R} a_{ij}$  and  $\sum_{j \in R} a_{ij} C_j$ . Equation (45) confirms that only one treated unit will be assigned to a single set



and Equation (46) confirms that a maximum of  $M$  control units can be assigned to a single set. Equation (47) says that the variable  $a_{ij}$  can take value 1 only if the value of parameter  $\text{dist}_{ij}$  is 1.

For minimizing the test statistic we need several extra equations. We first linearize the absolute value term  $\left| \sum_M \sum_{m=1}^M (n_{1,m-1,M} - \frac{m}{M+1} T_{m,M}) \right|$  by adding the following four additional constraints (50)-(53) with two new decision variables  $\psi_{m,M}$  and  $\omega_{m,M}$ ; where  $\psi_{m,M}$  is nonzero when  $n_{1,m-1,M} > \frac{m}{M+1} T_{m,M}$  and  $\omega_{m,M}$  is nonzero when  $n_{1,m-1,M} < \frac{m}{M+1} T_{m,M}$  :

$$\psi_{m,M} \geq n_{1,m-1,M} - \frac{m}{M+1} T_{m,M} \quad \forall m, M \quad (50)$$

$$\psi_{m,M} \geq 0 \quad \forall m, M \quad (51)$$

$$\omega_{m,M} \geq \frac{m}{M+1} T_{m,M} - n_{1,m-1,M} \quad \forall m, M \quad (52)$$

$$\omega_{m,M} \geq 0 \quad \forall m, M. \quad (53)$$

Thus  $\left| \sum_M \sum_{m=1}^M (n_{1,m-1,M} - \frac{m}{M+1} T_{m,M}) \right|$  now equals  $\psi_{m,M} + \omega_{m,M}$ .

As we are optimizing the value of  $\chi$ , we need to formulate the numerator, which is  $\left| \sum_M \sum_{m=1}^M (n_{1,m-1,M} - \frac{m}{M+1} T_{m,M}) \right| - 1/2$ . We add the following constraints (54)-(57) with two new decision variables  $\varphi_{m,M}$  and  $v_{m,M}$ , where  $\varphi_{m,M} > 0$  when the absolute value is greater than 1/2 and  $v_{m,M} > 0$  when the opposite holds. Finally, the resultant model for minimizing  $\chi$  is as follows:

$$\text{Minimize}_{\mathbf{a}} \quad \sum_M \sum_{m=1}^M (\varphi_{m,M} + v_{m,M})$$

subject to (23)-(53) and:

$$\varphi_{m,M} \geq \sum_M \sum_{m=1}^M (\psi_{m,M} + \omega_{m,M}) - 1/2 \quad \forall m, M \quad (54)$$

$$\varphi_{m,M} \geq 0 \quad \forall m, M \quad (55)$$

$$v_{m,M} \geq 1/2 - \sum_M \sum_{m=1}^M (\psi_{m,M} + \omega_{m,M}) \quad \forall m, M \quad (56)$$

$$v_{m,M} \geq 0 \quad \forall m, M. \quad (57)$$

In addition to the above constraints user may also choose to add covariate balance constraint. Similarly to the  $z$ -test algorithm, we will solve a series of ILP instances by tuning the

value of  $\phi$  to get a guaranteed optimal solution. We then get the optimal value of  $\chi$  by taking the optimal ratio of the objective value and corresponding  $\phi$  value, and square it to obtain the test statistic  $\chi^2$ . We note that the formulations we propose are not unique, there are other ways to encode the same optimization problems. We provided formulations that are linear and tend to work better in practice than some of the other formulations we tried.

In previous work, Hansen and Olsen Klopfer (2006), Hansen (2004) proposed methods for 1:M matches and applied them to evaluation of coaching for the SAT exam.

## 6. Case Studies

Our optimization models have been implemented in AMPL (Fourer et al. 2002), and solved with the solver CPLEX (ILOG 2007). The following four case studies demonstrate our proposed algorithms. All three datasets are publicly available for the purpose of reproducibility. The reported solution time with a X64-based PC with Intel(R) Core(TM) i7-4790 CPU running at 3.60 GHz with 16 GB memory and mip gap of 0.001 to solve a single instance is less than 1 second for all the tests, except  $\chi^2$  test, which was close to 1 minute.

### 6.1. Case Study 1

In our first case study, we used 2 years (2011-2012) of bike sharing data from the Capital Bike Sharing (CBS) system (see Fanaee-T and Gama 2014) from Washington DC comprised of 3,807,587 records over 731 days, from which we chose 247 treatment days and 463 control days according to the weather as follows: The control group consists of days with Weather 1: Clear, Few clouds, Partly cloudy, and Partly cloudy. The treatment group consists of days with Weather 2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist. The covariates for matching are as follows: Season (Spring; Summer; Fall; Winter), Year (2011; 2012), Workday (No; Yes), Temperature (maximum 41 degree Celsius), Humidity (maximum 100 percent), Wind speed (maximum 67). The outcome is the total number of rental bikes. We computed distance between days as follows:  $\text{dist}_{ij} = 1$  if covariates season, year and workday were the same, and the differences in temperature, humidity

and wind speed are less or equal to 2, 5 and 5, respectively for treated unit  $i$  and control unit  $j$ , 0 otherwise.

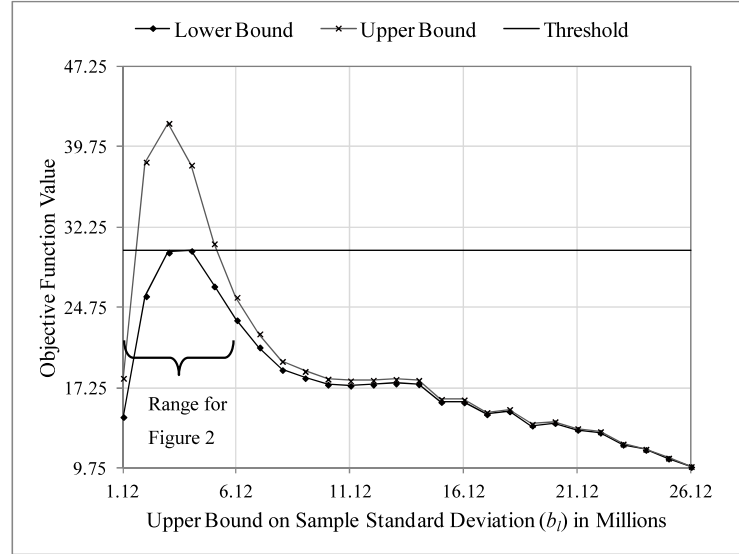
Figures 1 through 3 show the upper and lower bounds for the maximum objective function value for different  $b_i$  with  $n=30$ . These figures illustrate the meshes at different scales within the algorithm for computing the optimal solution for the maximization problem. Figures 4 through 6 illustrate the same phenomenon for minimization of the objective function. To produce Figure 7, we used similar procedure for producing Figures 1-6 at different values of  $n$ , namely  $n=30, 50, 70, 90, 110$ . The problem of finding pairs becomes infeasible above  $n = 110$ . Then each z-score was translated into a  $P$ -value  $1 - \phi(z)$  for both the upper and lower bounds for each  $n$ . The  $P$ -value for the upper bounds are very close to 1 and the  $P$ -value for the lower bounds are close to 0, illustrating that there is a lot of uncertainty associated with the choice of experimenter – one experimenter choosing 90 matched pairs can find a  $P$ -value of  $\sim 0$  and declare a statistically significant difference while another experimenter can find a  $P$ -value of  $\sim 1$  and declare the opposite. In this case it is truly unclear whether or not mist has an effect on the total number of rental bikes.

A note of caution: if the experimenter simply chose the  $P$ -value for the largest set of matched pairs, this would have deterministically yielded a  $P$ -value of almost 1 for 110 matched pairs. This is potentially very misleading – if we asked for slightly fewer pairs the result becomes completely unclear. Essentially, this casts doubt on the robustness of this particular choice used in standard practice.

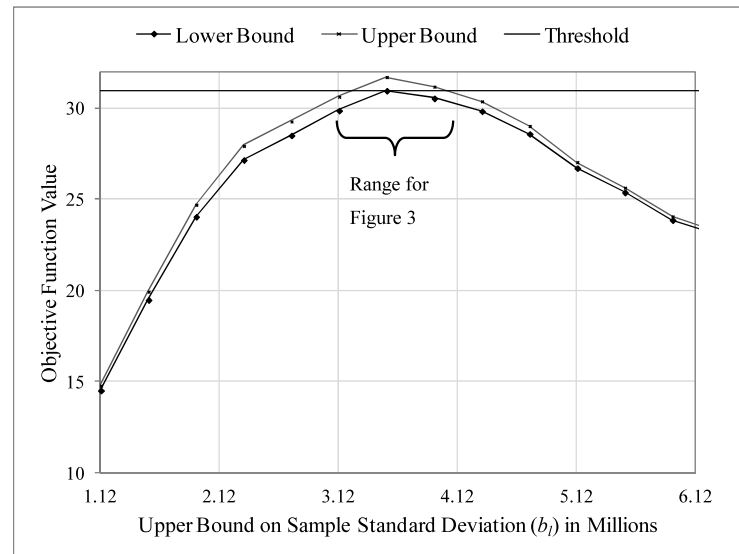
## 6.2. Case Study 2

In this case study we have used the data from a U.S. Department of Justice study regarding crime during the transition to adulthood, for youths raised in out-of-home care (see Courtney and Cusick 2010). Each observation represents a youth, and the outcome is whether he or she committed a violent crime over the 3 waves of the study.

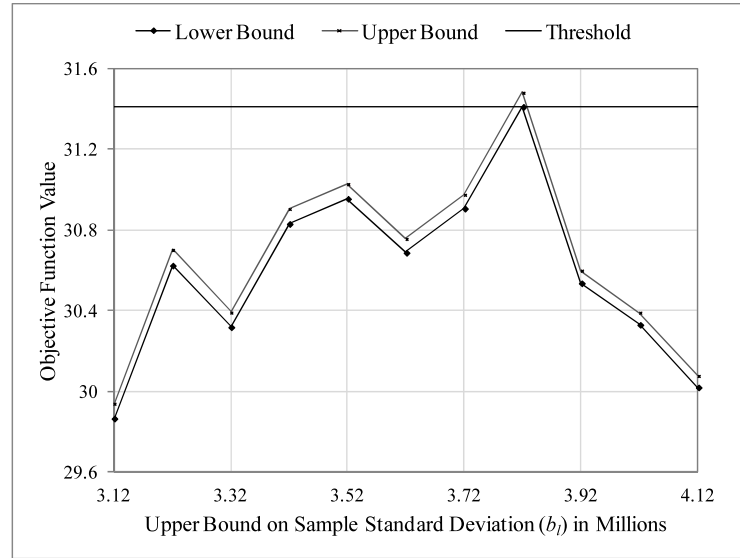
The (binary) covariates are as follows: hispanic, white, black, other race, alcohol or drug dependency, mental health diagnosis, history of neglect, entry over age of 12, entry age under 12, in school



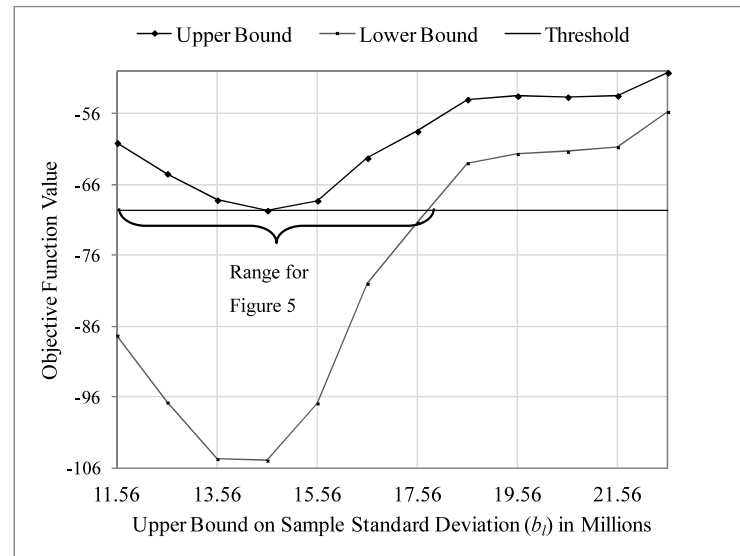
**Figure 1** Upper and lower bounds for maximum  $z$ -test objective function value over a range of  $b_l$  (Case Study 1:  $n=30$ ), illustrating the range for the finer mesh in the next step of the algorithm, which is found in Figure 2.



**Figure 2** Upper and lower bounds for maximum  $z$ -test objective function value over a range of  $b_l$  (Case Study 1:  $n=30$ ), illustrating the range for the finer mesh in the next step of the algorithm, which is found in Figure 3.

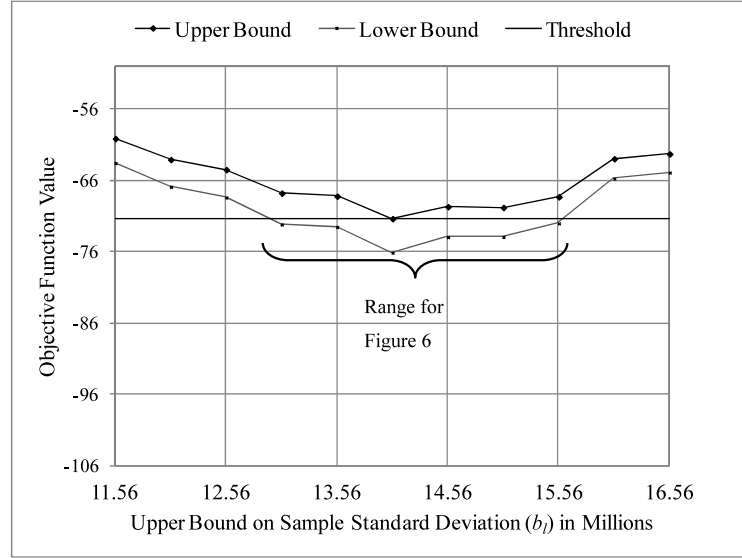


**Figure 3** Upper and lower bounds for maximum  $z$ -test objective function value over a range of  $b_l$  (Case Study 1:  $n=30$ ), at the finest mesh. The final value for the maximization problem is between 31.41 and 31.48.

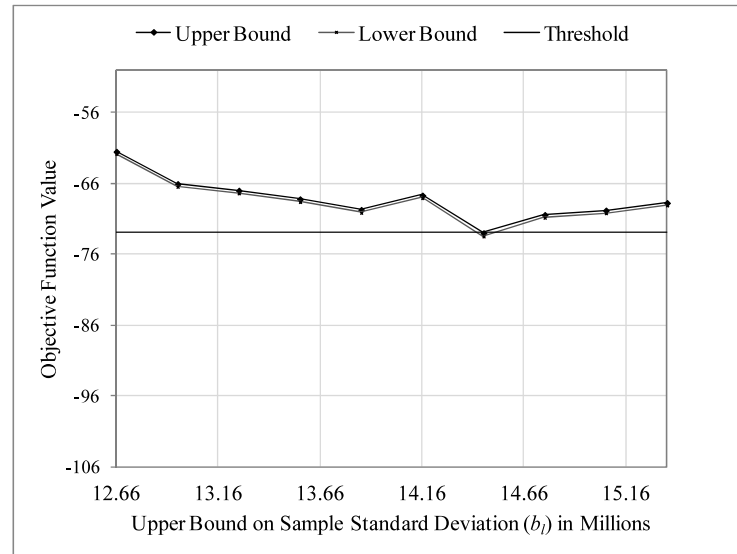


**Figure 4** Upper and lower bounds for minimum  $z$ -test objective function value over a range of  $b_l$  (Case Study 1:  $n=30$ ), illustrating the range for the finer mesh in the next step of the algorithm, which is found in Figure 5.

or employed, prior violent crime, any prior delinquency. The “treatment” variable is whether or not the individual is female; in particular we want to determine whether being female (controlling for race, criminal history, school/employment and relationship with parents) influences the probability

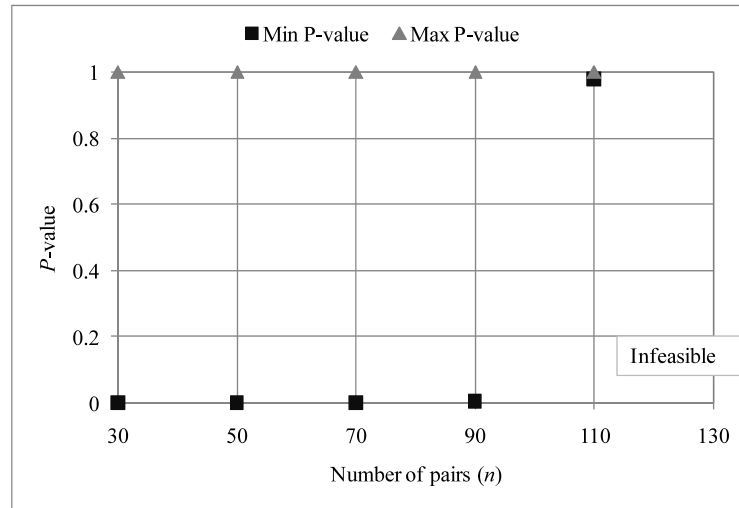


**Figure 5** Upper and lower bounds for minimum  $z$ -test objective function value over a range of  $b_l$  (Case Study 1:  $n=30$ ), illustrating the range for the finer mesh in the next step of the algorithm, which is found in Figure 6.

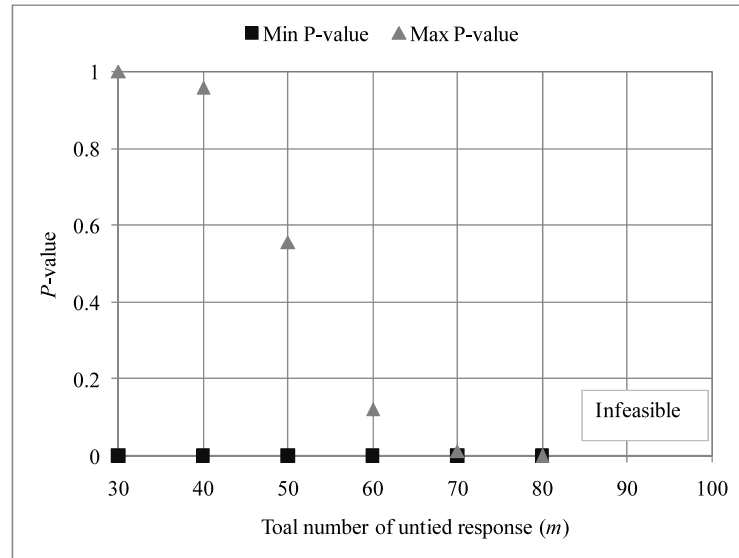


**Figure 6** Upper and lower bounds for minimum  $z$ -test objective function value over a range of  $b_l$  (Case Study 1:  $n=30$ ), at the finest mesh. The final value for the minimization problem is between -73.01 and -73.47.

of committing a violent crime. Here  $\text{dist}_{ij} = 1$  whenever all covariates of treated unit  $i$  are the same as those of the control unit  $j$ , 0 otherwise.

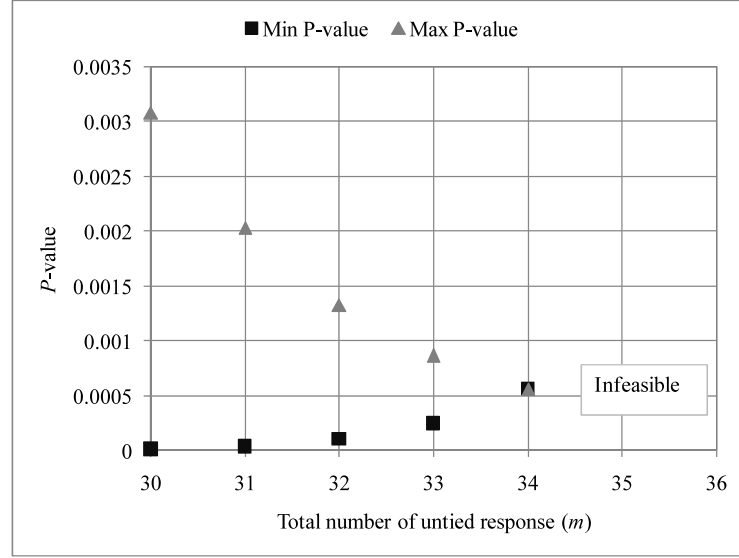


**Figure 7** Variation of  $z$ -test optimum  $P$ -values for different  $n$ . (Case Study 1)



**Figure 8** Variation of McNemar's test optimum  $P$ -value for different  $m$ . (Case Study 2)

Figure 7 is constructed in an analogous way to Figure 8 (using McNemar's test rather the  $z$ -test) showing the total number of discordant pairs along the x axis. Here, any matched pairs assignment would show a significant difference for the risks of violence between males and females. This difference becomes more pronounced as the number of pairs increases. Thus, the outcome is robust to the choice of experimenter.



**Figure 9** Variation of McNemar's test optimum  $P$ -value for different  $m$ . (Case Study 3)

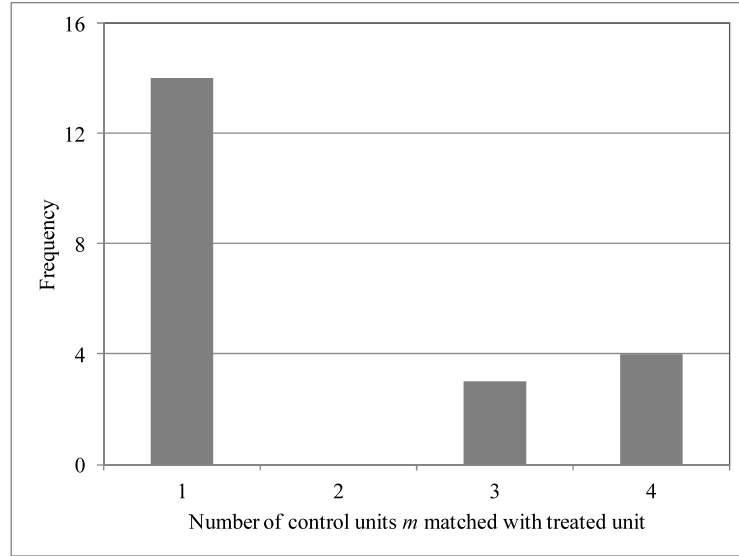
### 6.3. Case Study 3

In this case study we used GLOW (*Global Longitudinal study of Osteoporosis in Women*) data used in the study of Hosmer et al. (2013). Each row of data represents a patient and the outcome is whether or not the person developed a bone fracture. The treatment is smoking. The covariates are: age, weight, height, and BMI. Here  $\text{dist}_{ij} = 1$  whenever the difference between covariates of treated unit  $i$  and control unit  $j$  are all less than or equal to 6, and  $\text{dist}_{ij} = 0$  otherwise.

Figure 9 indicates robustly that smoking causes fracture, no matter which experimenter conducts McNemar's test.

We also ran the robust  $\chi^2$  test allowing up to  $M = 4$  control units per matched set. The minimum  $\chi^2$  value was 0, and the maximum was 5.2, where the solution for the maximum had several sets with  $m > 1$ . In particular, the count of pairs for the maximum solution is shown in Figure 10, illustrating that several treatment observations were matched with  $m = 4$  control observations, which gives them more weight in the solution. This unevenness of  $m$  could result naturally, depending on the density of treatment and control observations within the space  $x$ .



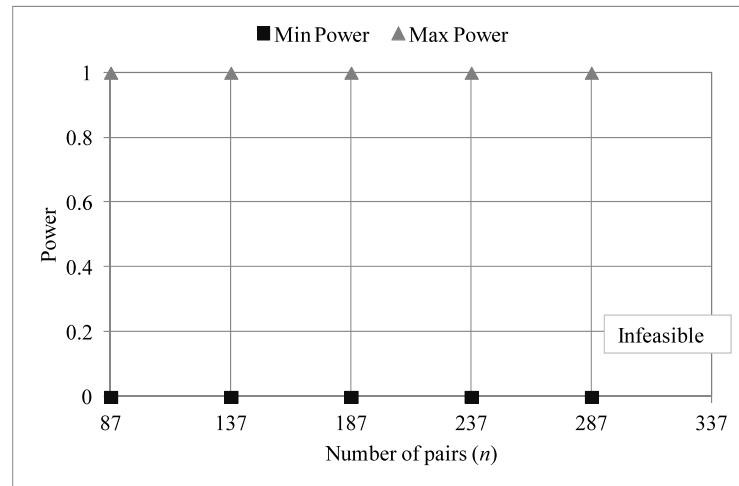


**Figure 10** Frequency of number of control units  $m$  matched with treated unit obtained by maximizing  $\chi^2$  test statistic. (Case Study 3)

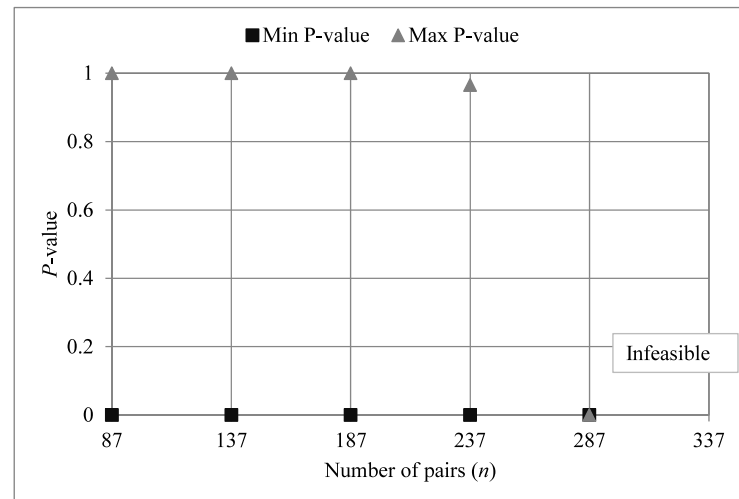
#### 6.4. Case Study 4

In this case study, we used training program evaluation data described in Dehejia and Wahba (1999), and Dehejia and Wahba (2002), which were drawn from Lalonde (1986). This data set contains 15,992 control units and 297 treatment units. The covariates for matching are as follows: age, education, Black (1 if black, 0 otherwise), Hispanic (1 if Hispanic, 0 otherwise), married (1 if married, 0 otherwise), nondegree (1 if no degree, 0 otherwise), and earnings in 1975. The outcome is the earnings in 1978. The treatment variable is whether an individual receives job training. We computed distance between units as follows:  $\text{dist}_{ij} = 1$  if covariates Black, Hispanic, married and nondegree were the same, and the differences in age, education and earnings in 1975 were less or equal to 5, 4 and 4000, respectively, for treated unit  $i$  and control unit  $j$ , 0 otherwise.

In the Figure 11, the  $P$ -value upper bounds are 1 and the  $P$ -value lower bounds are 0, illustrating that there is a lot of uncertainty associated with the choice of experimenter – one experimenter choosing 287 matched pairs can find a  $P$ -value of  $\sim 0$  and declare a statistically significant difference while another experimenter can find a  $P$ -value of  $\sim 1$  and declare the opposite. In this case it is truly unclear whether or not training has an effect on the earnings. To sanity check whether



**Figure 11** Variation of  $z$ -test optimum  $P$ -values for different  $n$ . (Case Study 4)



**Figure 12** Variation of  $z$ -test optimum  $P$ -values for different  $n$ . (Case Study 4 with additional random noise on the treatment outcome)

a reasonably sized effect would have been detected had one been present, we injected synthetic random noise (with normal distribution of mean  $\simeq$  \$10,000 and standard deviation  $\simeq$  \$100) on the treatment outcome, and the  $z$ -test (see Figure 12) robustly detects the treatment effect before the solutions become infeasible.

## 7. Conclusions

Believing hypothesis test results conducted from matched pairs studies on observational data can be dangerous. These studies typically ignore the uncertainty associated with the choice of experimenter, and in particular, how the experimenter chooses the matched pairs. In one of our case studies, we showed that it is possible to construct matched pairs so that the treatment seems to be effective with high significance, and yet another set of matched pairs exists where the estimated treatment effect is completely insignificant. We want to know that for *any* reasonable choice of experimenter who chooses the assignments, the result of the test would be the same.

In this work, we considered the most extreme sets of assignments that can be constructed, in a computationally efficient way involving integer linear programs. These are the most extreme hypothesis test results from reasonable matched pairs. Across the sciences, medical informatics, politics, and in other areas, if we not only knew that the test was significant for *some* experimenter, but that it was significant for *all* reasonable experimenters, then the result would be much more trustworthy, and relevant to policy makers.

In the future, we have been extending this work to nonparametric hypothesis tests. Another avenue for further research is to consider tests applying to evidence factors, where units are given different levels of treatment assignment (see Rosenbaum 2010b, 2011).

## Acknowledgments

The authors express their gratitude to the Natural Sciences and Engineering Research Council of Canada (NSERC) for their financial support of this research.

## References

- Breslow, N.E., N.E. Day. 1980. Statistical methods in cancer research. volume i - the analysis of case-control studies. *IARC scientific publications* **32** 5–338.
- Chen, D.-S., R.G. Batson, Y. Dang. 2011. *Applied Integer Programming: Modeling and Solution*. Wiley.
- Courtney, M.E., G.R. Cusick. 2010. Crime during the transition to adulthood: How youth fare as they leave out-of-home care in illinois, iowa, and wisconsin, 2002-2007. ICPSR27062-v1. Ann Arbor, MI: Inter-university Consortium for Political and Social Research.

- Dehejia, R.H., S. Wahba. 1999. Causal effects in nonexperimental studies: Reevaluating the evaluation of training programs. *Journal of the American Statistical Association* **94**(448) 1053–1062.
- Dehejia, R.H., S. Wahba. 2002. Propensity score-matching methods for nonexperimental causal studies. *Review of Economics and Statistics* **84**(1) 151–161.
- Fanaee-T, Hadi, Joao Gama. 2014. Event labeling combining ensemble detectors and background knowledge. *Progress in Artificial Intelligence* **2**(2-3) 113–127.
- Fourer, R., D.M. Gay, B.W. Kernighan. 2002. *Ampl: A modeling language for mathematical programming*. Duxbury Press, Cole Publishing Co.
- Hansen, Ben B. 2004. Full matching in an observational study of coaching for the SAT. *Journal of the American Statistical Association* **99** 609–618.
- Hansen, Ben B., S. Olsen Klopfer. 2006. Optimal full matching and related designs via network flows. *Journal of Computational and Graphical Statistics* **15**(3) 609–627.
- Ho, Daniel E., Kosuke Imai, Gary King, Elizabeth A. Stuart. 2007. Matching as nonparametric preprocessing for reducing model dependence in parametric causal inference. *Political Analysis* **15**(3 (Summer)) 199–236.
- Holland, Paul W. 1986. Statistics and causal inference. *Journal of the American Statistical Association* **81**(396) 945–960.
- Hosmer, D.W., S. Lemeshow, R.X. Sturdivant. 2013. *Applied Logistic Regression: Third Edition..* John Wiley and Sons Inc.
- ILOG. 2007. Cplex 11.0 user’s manual. ILOG, Inc.
- Lalonde, R. 1986. Evaluating the econometric evaluations of training programs. *American Economic Review* **76**(4) 604–620.
- Morgan, Stephen L., Christopher Winship. 2007. *Counterfactuals and Causal Inference: Methods and Principles for Social Research*. Cambridge University Press.
- Noor-E-Alam, M., A. Mah, J. Doucette. 2012. Integer linear programming models for grid-based light post location problem. *European Journal of Operational Research* **222**(1) 17–30.

- 
- Rosenbaum, Paul R. 2012. Optimal matching of an optimally chosen subset in observational studies. *Journal of Computational and Graphical Statistics* **21** 57–71.
- Rosenbaum, P.R. 1989. Optimal matching for observational studies. *Journal of the American Statistical Association* **84** 1024–1032.
- Rosenbaum, P.R. 2010a. *Design of Observational Studies*. Springer, New York.
- Rosenbaum, P.R. 2010b. Evidence factors in observational studies. *Biometrika* **97**(2) 333–345.
- Rosenbaum, P.R. 2011. Some approximate evidence factors in observational studies. *Journal of the American Statistical Association* **106**(493) 285–295.
- Rubin, D.B. 2007. The design versus the analysis of observational studies for causal effects: Parallels with the design of randomized trials. *Statistics in Medicine* **26**(1) 20–36.
- Rubin, D.B. 2008. For objective causal inference, design trumps analysis. *Annals of Applied Statistics* **2**(3) 808–840.
- Rubin, Donald B. 1974. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology* **5**(66) 688–701.
- Tamhane, A.C., D.D. Dunlop. 2000. *Statistics and Data Analysis..* Prentice Hall, New Jersey.
- Winston, W.L., M. Venkataramanan. 2003. *Introduction to Mathematical Programming, (4th ed.)*. Thomson (Chapter 9).
- Wolsey, L.A. 1998. *Integer Programming*. Wiley-Interscience, Series in Discrete Mathematics and Optimization, Toronto.
- Zubizarreta, J.R. 2012. Using mixed integer programming for matching in an observational study of kidney failure after surgery. *Journal of the American Statistical Association* **107**(500) 1360–1371.
- Zubizarreta, J.R., R.D. Paredes, P.R. Rosenbaum. 2014. Matching for balance, pairing for heterogeneity in an observational study of the effectiveness of for-profit and not-for-profit high schools in chile. *The Annals of Applied Statistics* **8**(1) 204–231.
- Zubizarreta, J.R., D.S. Small, N.K. Goyal, S. Lorch, P.R. Rosenbaum. 2013. Stronger instruments via integer programming in an observational study of late preterm birth outcomes. *Annals of Applied Statistics* **7**(1) 25–50.

## 8. Supplement A

### 8.1. Integer Linear Programming Basics

ILP techniques have become practical for many large-scale problems over the past decade, due to a combination of increased processor speeds and better ILP solvers. Any type of logical condition can be encoded as linear constraints in an ILP formulation with binary or integer variables. Consider two binary variables  $x \in \{0, 1\}$  and  $y \in \{0, 1\}$ . The logical condition “if  $y = 0$  then  $x = 0$ ” can be simply encoded as

$$x \leq y.$$

Note that this condition imposes no condition on  $x$  when  $y = 1$ . Translating if-then constraints into linear constraints can sometimes be more complicated; suppose, we would like to encode the logical condition that if a function  $f(w)$  is greater than 0, then another function  $g(w)$  is greater or equal to 0. We can use the following two linear equations to do this, where  $\theta$  is a binary variable and  $M$  is a positive number that is larger than the maximum values of both  $f$  and  $g$ :

$$-g(w) \leq M\theta$$

$$f(w) \leq M(1 - \theta).$$

In order for  $f(w)$  to be positive, then  $\theta$  must be 0, in which case,  $g(w)$  is then restricted to be positive. If  $f(w)$  is negative,  $\theta$  must be 0, in which case no restriction is placed on the sign of  $g(w)$ . (See for instance the textbook of Winston and Venkataramanan (2003), for more examples of if-then constraints).

ILP can capture other types of logical statements as well. Suppose we would like to incorporate a restriction such that the integer variable  $S_i$  takes a value of  $K$  only if  $i = t$ , and 0 otherwise. The following four if-then constraints can be used to express this statement, where  $\lambda_1$  and  $\lambda_2$  are binary variables:

$$\lambda_1 = 1 \text{ if } i + 1 > t$$

$$\lambda_2 = 1 \text{ if } t + 1 > i$$

$$S_i = k \text{ if } \lambda_1 + \lambda_2 > 1$$

$$S_i = 0 \text{ if } \lambda_1 + \lambda_2 < 2.$$

Each of these if-then constraints (4)-(7) can be converted to a set of equivalent linear equations, similar to what we described above. (See also Noor-E-Alam et al. (2012) and Winston and Venkataramanan (2003)).

There is no known polynomial-time algorithm for solving ILP problems as they are generally NP-hard, but they can be solved in practice by a number of well-known techniques (Wolsey (1998)). The LP relaxation of an ILP provides bounds on the optimal solution, where the LP relaxation of an ILP is where the integer constraints are relaxed and the variables are allowed to take non-integer (real) values. For instance, if we are solving a maximization problem, the solution of the LP relaxation can serve as an upper bound, since it solves a problem with a larger feasible region, and thus attains a value at least as high as that of the more restricted integer program. ILP solvers use branch-and-bound or cutting plane algorithms combined with other heuristics, and are useful for cases where the optimal integer optimal solution is not attained by the LP relaxation. The branch-and-bound algorithms often use LP relaxation and semi-relaxed problems as subroutines to obtain upper bounds and lower bounds, in order to determine how to traverse the branch-and-bound search tree (Chen et al. 2011, Wolsey 1998). The most popular ILP solvers such as CPLEX, Gurobi and MINTO each have different versions of branch-and-bound techniques with cutting plane algorithms and problem-specific heuristics.