

**Utilizing the Zero-One Linear Programming Constraints to Draw Multiple Sets of Matched Samples from a Non-Treatment Population as Control Groups for the Quasi-Experimental Design**

**By**

**Yuan H. Li, Yu, N. Yang & Leroy J. Tompkins**

**Prince George's County Public Schools, Maryland**

**Shahpar Modarresi**

**Montgomery County Public Schools, Maryland**

**Paper was presented at the annual meeting of the American Educational Research Association, Montreal, Canada, April, 2005.**

# Utilizing the Zero-One Linear Programming Constraints to Draw Multiple Sets of Matched Samples from a Non-Treatment Population as Control Groups for the Quasi-Experimental Design

Abstract: The statistical technique, *Zero-One Linear Programming*, that has successfully been used to create multiple tests with similar characteristics (e.g., item difficulties, test information and test specifications) in the area of educational measurement, was deemed to be a suitable method for creating multiple sets of matched samples to be used as control groups in the quasi-experimental design of *non-randomized comparison group pretest-posttest*. Compared to the existing propensity-score matching method, this method does not require any statistical models and assumptions and can handle the covariate of the pretest score more appropriately.

If the measurement error of the pretest-score mean of the treatment group is ignored, this method will generate a *unique* matched sample once the criteria for attempting to create two similar groups are determined. Otherwise, multiple sets of similar matched samples can be generated and the performance of the treatment group can be compared with each of the multiple matched samples using an appropriate statistical analysis. Afterwards, the mean of the effect size measure, taking the average of the effect size across replicated comparisons, can then be used to assess the efficacy of any program. This enhances our confidence level to decide whether a program is effective or not, compared to the finding resulting from a single comparison.

A description of *Zero-One Linear Programming* and its application to create a matched sample or multiple sets of matched samples is introduced in this paper.

Key Words: Linear Programming, Matched Sample, Quasi-Experimental Design Optimization, Experimental Design, Program Evaluation

# I Introduction

## A. Background of Quasi-experimental Design

In an experimental design, random assignment is an ideal sampling method to create experimental and control groups when a group of subjects is available. The subjects of the experimental group will receive a treatment; whereas, no specific treatment will be given to the subjects of the control group. The procedure of random assignment becomes a powerful technique for controlling all known and “unknown” extraneous variables because it makes both groups very similar at the beginning of an experiment, especially in cases where the sample size is large. Unfortunately, this method has often encountered implementation obstacles in the evaluation of educational programs because student enrollment in a specific program is not random, and as such, cannot be completely manipulated as can be done with random assignment in most instances. Accordingly, the quasi-experimental design, defined as an experiment without randomized assignment but involving the manipulation of independent variables (Isaac & Michael, 1995; Shadish, Cook & Campbell, 2002), becomes one of the alternatives used to determine a program effect.

The *non-randomized comparison group pretest-posttest design* (illustrated in Figure 1, refer to Shadish, Cook & Campbell [2002], p. 136) is the most appropriate evaluation design in assessing the efficacy of any program among the quasi-experimental designs. For this design, random assignment is not conducted and subjects in both the quasi-experimental and the control groups will take both the pretest and the posttest. Like a true experimental design, the subjects in the control group will not receive any specific treatment, but their counterparts in the quasi-experimental group(s) will receive program treatment(s). Here, the number of groups under the quasi-experimental label could be single (e.g., only one program) or multiple (e.g., several programs to be evaluated simultaneously).

Without random assignment in the design introduced above, the impact of undetected nuisance variables on the outcome variable might not be ruled out. In order to better assess the efficacy of a program, a purely statistical modeling (e.g., analysis of covariance, ANCOVA, Kirk, [1995]; hierarchical linear modeling, HLM, Bryk & Raudenbush, [1992]) may be used to tackle this issue. Without using the matched sample procedure, the statistical modeling is primarily used to account for student’s characteristic differences (e.g., sex, race, pretest scores, etc.) or school context differences (e.g., percent of minority students, percent of poverty students, etc.) between treated and non-treatment groups. However, as pointed out by Rubin, Stuart and Zanutto (2004), comparing results obtained from treated (e.g., magnet programs) and whole control groups (e.g., non-magnet population) with very different distributions of background covariates will heavily rely on untestable modeling (e.g., ANCOVA) assumptions and extreme extrapolation. As such, reliable causal inferences may not be drawn. For example, Rubin et. al. (2004) further illustrated that the values of “percent minority” and “percent in poverty” may differ widely at some schools, this situation will cause the estimated school effects that have been adjusted for such covariates using models be extremely sensitive to these statistical modeling assumptions (e.g., parallel slopes). If the assumptions are seriously violated the distributions of background variables among subgroups are different to some extent, the estimated program effect, as a result of using extreme extrapolation, will be seriously misleading.

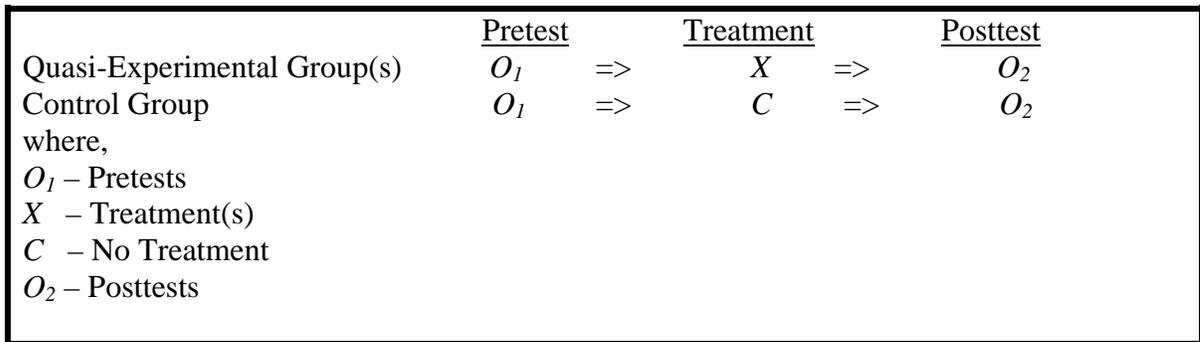


Figure 1: The Non-Randomized Comparison Group Pretest-Posttest Design

### B. Issues Associated with Creating a Matched Sample

The goal of obtaining a *Matched Sample as a Control Group* is to create the conditions similar to a randomized experiment as closely as possible. The treated and control groups are matched without using the observed outcome variable (or posttest), thus preventing us from “intentionally” manipulating a matched sample to obtain a desired result and also protecting from such claims by researchers. The ability of a matched sample procedure to reveal the extent to which treated and matched groups have similar types of students in similar educational settings is “an important diagnostic tool to identify whether the data can support [possible] causal comparisons between these two groups” (Rubin et. al., 2004).

Instead of totally utilizing statistical modeling in assessing program effects, it would be preferable to use an appropriate matched sample to be compared with the treatment group before any statistical modeling is performed. However, creating an appropriate matched sample for each program is a major challenge. Diverse methods can be used to accomplish this objective (for literature review, see Shadish et. al., 2002). Among them, using the propensity score (Rosenbaum & Rubin, 1983, 1984, 1985; Rosenbaum, 1995) as a criterion for selecting students as members of a matched sample is one of the promising approaches to address this issue. A propensity score is an estimated probability of a given individual belonging to a treatment group given the observed background characteristics (or covariates) of that individual. This propensity score reduces the entire collection of background characteristics to a single composite index value so that a matched sample will be selected using this single index, instead of directly matching multiple background variables.

Nevertheless, the value of the propensity score is dependent on the selections of statistical models (e.g., whether or not including the interaction, and/or nonlinear terms on the logistic regression models). Also, if the assumptions made for the statistical model (e.g., logistic regression) are not met, and/or if the sample size used for the model is not large enough, using those propensity scores as a criterion for selecting a matched sample might not be as meaningful as researchers anticipated. Furthermore, the weighting for each covariate, that is then used for computing the propensity score, depends totally on the degree of each covariate’s relation to the treatment assignment (received or not received treatment). This procedure is not appropriate for the *non-randomized comparison group pretest-posttest design*, in which the pretest score is usually highly correlated with the outcome measure rather than the treatment assignment. This scenario will cause the



□ is the relationship function, which could be  $\leq$ ,  $=$ , or  $\geq$ . The equal symbol of ‘=’ is used here.

More specifically, in Equation 1, the members in the population are indexed by  $i=1, \dots, n$  and the values in the variable  $x_i$  are parameters that will be estimated. For *zero-one linear programming*, the  $x$  values are constrained to be either one or zero as indicated in Equation 3 to identify whether the members are selected or not for the matched group.

Equation 2 introduced above can be presented by a matrix expression—Equation 4 (shown below) in which the vector of  $\underline{x}$  will be resolved by not only maximizing (or minimizing) the linear function of Equation 1, but also imposing the constraint of  $x$  values of either one or zero. The matrix  $\mathbf{A}$  and the vector  $\underline{b}$  in Equation 4 are created from  $A_{m \times n}$  and  $b_m$  coefficients, respectively. The way of preparing both matrix  $\mathbf{A}$  and the vector  $\underline{b}$  depends on the nature of the problem we attempt to resolve. It is noted that the following descriptions in illustrating how to prepare both matrix  $\mathbf{A}$  and the vector  $\underline{b}$  for the solution of *zero-one linear programming* only fit the problem presented in this paper. Readers might refer to other references (e.g., Theunissen, 1985, 1986) for better understanding this issue.

$$\mathbf{A} \cdot \underline{x} = \underline{b} \tag{4}$$

It is noted that if multiple pretest scores are available, a composite score obtained from those pretests is more appropriate to be entered into Equation 1. The choice of types of composite score can be dependent on the nature of those pretest scores themselves.

### B. Example of Using Zero-One Linear Programming

In the present example, suppose two key student demographic variables (e.g., gender and poverty) are considered for matching. Under this circumstance, there are four types of students as shown in Table 1 – male/poverty, male/non-poverty, female/poverty, and female/non-poverty. It is further assumed that there are 10 and 20 students in the magnet and non-magnet (or comprehensive) programs, respectively, and the frequencies of each type of student are also shown in Table 1. Ten non-magnet students will be drawn from across the four student subgroups to correspond with the 10 magnet students. The number of non-magnet students drawn from each subcategory will be identical to the number of magnet students in the respective subcategory. At the same time, the average pretest scores between two groups are expected to be as close as possible.

Table 1. An Example of Data Regarding the Frequency of Four Types of Students

Type of Students	Magnet Students Frequency	Non-Magnet Students Frequency
1. male/poverty	2	4
2. male/non-poverty	3	6
3. female/poverty	1	2
4. female/non-poverty	4	8
<b>Subtotal</b>	<b>10</b>	<b>20</b>

Under the sampling scenario cited above, the  $\mathbf{A}$  matrix and  $\underline{\mathbf{b}}$  vector in Equation 2 or 4 will be created as shown below before seeking the “Zero-One” solution of the vector parameter  $\underline{\mathbf{x}}$ .

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

$$\underline{\mathbf{b}} = \begin{bmatrix} 2 \\ 3 \\ 1 \\ 4 \end{bmatrix}$$

Regarding the  $\mathbf{A}$  matrix, the number of columns in the  $\mathbf{A}$  matrix should equal the number of non-magnet students in the pool. In addition, each row in the  $\mathbf{A}$  matrix together with the corresponding row in the  $\underline{\mathbf{b}}$  vector expresses a single constraint. The first constraint is expressed in the first row in the  $\mathbf{A}$  matrix together with the first row in the  $\underline{\mathbf{b}}$  vector. The four series of “1” connotes that the first four of the 20 students are Type 1 students, and the rest of the sixteen series of “0” connotes that they are not Type 1 students. Further, the condition of two members in Type 1 students to be picked as part of a matched sample is specified as “2” in the first row in the  $\underline{\mathbf{b}}$  vector.

The second constraint is expressed in the second row in the  $\mathbf{A}$  matrix together with the second row in the  $\underline{\mathbf{b}}$  vector. The six series of “1” connotes that they are Type 2 students and the rest of the fourteen of “0” connotes that they are not Type 2 students. Further, the condition of three members in Type 2 students to be picked as part of a matched sample is specified as “3” in the second row in  $\underline{\mathbf{b}}$  vector. Using the same logic, the third and fourth constraints are specified in the  $\mathbf{A}$  matrix and the  $\underline{\mathbf{b}}$  vector.

After the  $\mathbf{A}$  matrix and the  $\underline{\mathbf{b}}$  vector are set up and both are then inserted into Equation 4 ( $\mathbf{A} \cdot \underline{\mathbf{x}} = \underline{\mathbf{b}}$ ), a specific mathematical function is formed and shown in Equation 5. The solution of *Zero/One* values of  $x_i$  parameters in this function will be found on the condition that the targeted function of Equation 1 is minimized in this case.

$$\begin{bmatrix} 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \\ x_{10} \\ x_{11} \\ x_{12} \\ x_{13} \\ x_{14} \\ x_{15} \\ x_{16} \\ x_{17} \\ x_{18} \\ x_{19} \\ x_{20} \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \\ 1 \\ 4 \end{bmatrix} \tag{5}$$

As already indicated, a statistical technique, *Zero-One Linear Programming*, is used to seek the solution of the vector of  $\underline{x}$ . Afterwards, a matched sample will be created using the information of the  $x$  *zero/one* values. This matched sample will have the same demographic characteristics (e.g., poverty status, gender) as the magnet student population because the constraints formalized in Equations 2 and 3 are imposed. Furthermore, its average pretest score is close to the average pretest score of the magnet population because the linear function in Equation 1 is minimized.

*It is important to note that “the same demographic characteristics” in those selected matching variables (e.g., poverty and gender variables) not only means that the number of students in each selected variable itself is identical between two groups, but also means that the distribution of students in the combination (e.g., four types of students cited above) of those selected matching variables (e.g., poverty and gender) is identical. The latter feature has been demonstrated in the above example, but is too complicated to be done using previously existing matching procedures.*

For the example introduced above, only the category variables (e.g., gender, race) were selected to create a series of constraints formulized in Equation 2. Other continuous variables (e.g., age) are also suitable to be included (refer to Theunissen, 1985, 1986).

### **III. A First Study Using the *Zero-One Linear Programming* Approach to Create a Matched Sample**

An investigation of the effects of magnet school programs on the reading and mathematics performance of students in a school system was conducted by Yang, Li, Modarresi and Tompkins (2003). The major objective of this summative evaluation was to compare the academic performance of students from each of the magnet programs, to the performance of a matched sample of their non-magnet peers in the reading and mathematics content areas. The *non-randomized comparison group pretest-posttest design* has been used for this evaluation. Refer to Figure 1, the quasi-experimental group was a group of magnet program students, and the control group was a matched sample that was drawn from the population of non-magnet students using the ***Zero-One Linear Programming*** technique introduced above. The magnet program group received the magnet program treatment, while the non-magnet group did not.

The Comprehensive Test of Basic Skills (CTBS, CTB/McGraw-Hill, [1997]) reading and mathematics tests administered in 2001 were used as pretests for both groups of students. The 2003 Maryland School Assessment (MSA) reading and mathematics assessments were used as posttests. Measuring the magnet program treatment effect was of primary interest in the study cited above. The posttest score difference between two groups might be used for this purpose; however, the pre-existing difference between the two groups (e.g., initial abilities and demographic differences) was not accounted for by only observing the simple posttest score difference. Most researchers often suggest that, if possible, both statistical and matching controls should be simultaneously employed in order to better adjust for those pre-existing differences at the beginning of the experiment. Such a principle was fully applied on that study, in which a matched sample was created for each magnet program before the ANCOVA was used to adjust for the small difference in the pretest score between the magnet and respective matched groups. The process of creating a matched sample (e.g., for the Academic Center Magnet Program) is described below.

#### **A. Tabulate the frequencies of various types of students**

As seen in Table 2, there were 468 and 6,435 fifth graders in the Academic Center Magnet Program and comprehensive (or Non-Magnet) programs, respectively. Students were grouped by combinations of race, gender, and poverty status, i.e., 20 types of students were classified and listed in the first column of Table 2. The frequencies of those 20 types of students are also shown in Table 2 for both Magnet and Non-Magnet groups.

Table 2.

Frequencies for the Academic Center Magnet Program and Non-Magnet Students

Types of Students	Magnet Program Frequency	Non-Magnet Students Frequency
1. American Indian, male, non-poverty	1	5
2. American Indian, male, poverty	0	7
3. American Indian, female, non-poverty	3	9
4. American Indian, female, poverty	0	8
5. Asian, male, non-poverty	3	53
6. Asian, male, poverty	3	34
7. Asian, female, non-poverty	7	55
8. Asian, female, poverty	2	32
9. African American, male, non-poverty	92	1,098
10. African American, male, poverty	80	1,399
11. African American, female, non-poverty	100	1,121
12. African American, female, poverty	103	1,548
13. White, male, non-poverty	19	215
14. White, male, poverty	0	46
15. White, female, non-poverty	19	219
16. White, female, poverty	4	44
17. Hispanic, male, non-poverty	9	35
18. Hispanic, male, poverty	15	210
19. Hispanic, female, non-poverty	2	46
20. Hispanic, female, poverty	6	251
<b>Total</b>	<b>468</b>	<b>6,435</b>

**B. Choose the test score to be minimized**

Since a matched sample is to be drawn and then applied to either reading or mathematics performance evaluation for this magnet program, the average pretest score of both CTBS reading and mathematics T scores was used to be minimized in the context of *Zero-One Linear Programming*. The T score equals  $(50 + 10 \text{ times } z)$ , where  $z$  is the standard score of the reading or mathematics “scale” score. Averaging T scores in reading and mathematics, instead of averaging their scale scores, is done to ensure that the weighting in both content areas is equal when both scores were added up together and then were averaged.

**C. Utilize the Zero-One Linear Programming**

Once each student’s average pretest score and the distribution of various types of students for the Academic Center are available, the linear function (presented in Equation 1), matrix  $A$  and the vector  $b$  can be created. The *zero-one linear programming* then used them to seek the solution of the  $x$  vector indicated in Equation 1. The  $x$  vector was then used to identify which students were chosen to be part of the matched sample from among 6,435 non-magnet students. The frequency distributions and average pretest scores for the Academic Center and its matched sample are provided in Table 3, where it shows that the distributions of various types of students between the magnet program and matched groups

are identical. Furthermore, their average pretest scores are almost the same. No statistically significant difference was found in the average pretest scores between the two groups.

Table 3.  
Final Results: Frequency Distributions and Average Pretest Scores for the Academic Center and Its Matched Sample

Types of Students	Magnet Program Frequency	Matched Sample Frequency
1. American Indian, male, non-poverty	1	1
2. American Indian, male, poverty	0	0
3. American Indian, female, non-poverty	3	3
4. American Indian, female, poverty	0	0
5. Asian, male, non-poverty	3	3
6. Asian, male, poverty	3	3
7. Asian, female, non-poverty	7	7
8. Asian, female, poverty	2	2
9. African American, male, non-poverty	92	92
10. African American, male, poverty	80	80
11. African American, female, non-poverty	100	100
12. African American, female, poverty	103	103
13. White, male, non-poverty	19	19
14. White, male, poverty	0	0
15. White, female, non-poverty	19	19
16. White, female, poverty	4	4
17. Hispanic, male, non-poverty	9	9
18. Hispanic, male, poverty	15	15
19. Hispanic, female, non-poverty	2	2
20. Hispanic, female, poverty	6	6
<b>Total N</b>	<b>468</b>	<b>468</b>
<b>Average Pretest Score</b>	<b>50.53</b>	<b>50.52</b>
<b>Difference of Pretest Score</b>	<b>0.01*</b>	

\* P = .986

#### IV. Individually-based Matching

The matched sample generated by the above steps did not specifically identify the respective matched member given a treatment group (e.g., the magnet-program) member. As already indicated, this group-based matching procedure has been incorporated into the magnet program study (Yang, Li, Modarresi and Tompkins, 2003). In some instances, a matched sample generated by an individual matching procedure is preferred. For example, when the program effect is analyzed by the whole group and is then required to be analyzed by the disaggregated subgroups, an individual matching procedure makes it possible that each subgroup has its own matched sample to be compared with.

This section further discusses the steps on how to modify current matching procedure to serve the purpose mentioned above. There are several possible solutions. The algorithm described below is one of the promising procedures.

- (1). Start with the first individual member from the magnet group.
- (2). Utilize the current group matching procedure to draw a matched sample.
- (3). Find a member from the matched sample generated by Step 2 with the following conditions, a) Have the same type of member as this individual member indicated in Step 1, b) Have the closest pretest score to this individual member indicated in Step 1. Once this member is found, he/she will be the respective member of this individual member indicated in Step 1.
- (4). Replace any members (note: those drawn by Step 2) that have not been previously selected by Step 3 in the non-magnet population pool.
- (5). Add an additional constraint to the constraints that have been imposed in the *zero/one* linear model to ensure that the member being recently selected in Step 3 “must” be included in the next new-drawn matched sample. For matrix expression in the above example, if the second member from the non-magnet population is selected to be matched with the first member of the magnet group, the **A** matrix and **b** vector should be updated in the following way:

$$\mathbf{A} = \begin{bmatrix} 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\underline{\mathbf{b}} = \begin{bmatrix} 2 \\ 3 \\ 1 \\ 4 \\ 1 \end{bmatrix}$$

The last row in **A** matrix combined with the last row in **b** vector indicates that the second member is certain to be one member in the next new-drawn matched sample.

- (6). Repeat the steps 1-5 until all members in the magnet group have found their own respective members from the non-magnet population.

The above individually-based matching procedure begins with drawing a matched sample that meets all the desirable constraints and has a minimum pretest-score difference between both matched and respective treatment groups. A member who meets the criteria indicated in Step 3 is then selected from this matched sample, instead of directly from the full non-treatment population. Each of the next serial matched samples includes all previously selected members and consequently the last matched sample is an actual

matched sample which always meets all constraints. Also, each individual member from the matched sample corresponds to a member from the treatment (or magnet) group. The group-based matching procedure introduced previously will often generate a matched sample with too little variability of the pretest scores, compared to the variability occurred in the treatment group. The individual matching procedure introduced in this section will alleviate this problem to some extent. The application of this matching procedure can be found on two program evaluation studies (e.g., Modarresi, Yang, Bulgakov-Cooke, & Li, 2004; Yang, Li, Modarresi & Tompkins, 2004). The other procedure to be introduced next will make not only the means of the pretest score but also their respective variances between two treatment and non-treatment groups very similar.

## V. Multiple-Matched-Sample Procedure

### A. Matching Procedure with the Involvement of Measurement Error

As indicated in the section of (group or individual) matching control, a unique matched sample will be generated once the criteria used for the matching procedure is determined. This is especially true if we assume the average pretest score of a treatment group (e.g., magnet program students) is a true score, not contaminated with any measurement error. For large sample sizes, this assumption should be appropriate. However, to increase the confidence level of seeking an appropriate matched sample as similar to the treatment group as we could obtain, such no-measurement-error is not necessary to be presumed by allowing the pretest-score mean to be contaminated with a “reasonable” measurement error. Equation 6 presented below will help us comprehend this concept.

$$\text{Minimize } \sum_{i=1}^n [\text{ABS}(P_i - (M + E))] x_i \quad (6)$$

The components in Equation 6 are the same as those found in Equation 1, except the additional component of measurement error, E. The value of E can be randomly generated from the normal distribution,  $N(0, SE^2)$ , where SE represents the standard error of the mean of pretest scores for the treatment group. Specifically,

$$SE^2 = \frac{S^2}{N} \quad (7)$$

Where

N is the sample size of treatment group,

$S^2$  is sample variance of pretest scores for the treatment group.

In reality, the matching procedure introduced in this section may still generate a matched sample whose pretest-score variance is still way off (too small) to the one found in the treatment group when the sample size of the treatment group, N, is too large. On the other hand, this matching procedure may generate a matched sample whose pretest-score mean is not very close to the one found in the treatment group when the sample size is too small. This issue can be manually resolved by: first conducting several trials for various sizes of N; second, finding an appropriate N that will produce a matched sample whose

mean as well as variance of the pretest score is very similar to the treatment group. Of course, another comprehensive approach can also be used for resolving this issue: first, decide how close the mean and variance of the pretest scores the two groups should be; second, repeatedly conduct the matching procedure until the criteria we set has been achieved. The iterative procedure that was often used in computer language (e.g., loops) can deal with this comprehensive approach very efficiently.

## **B. Multiple Sets of Matched Samples**

By allowing the addition of measurement error into the mean score of the pretest for the treatment group during the process of matching procedures, a matched sample will be generated. Afterwards, every member of the non-treatment group is returned to the dataset after sampling. Another matched sample will be generated given a different value of measurement error. Again, every member of the non-treatment group should be returned to the dataset after sampling. After repeating the matched procedure again and again, multiple matched samples will be created. It is noted that many members from the population of the non-treatment group could appear multiple times in different sets of matched samples because the same constraints have been repeatedly imposed into the matching procedure.

The multiple-matched sample procedure creates a condition that the treatment group has multiple matched samples to compare with. As performed in the one-matched-sample approach, an effect size measure (for the discussions of the features of this measure, refer to Thompson [2002]) can be performed for each analysis in each comparison. If 1000 matched samples are used as in this evaluation, the mean as well as the distribution of the effect size measure, across 1000 replicated comparisons, can be used to assess the efficacy of any program. This enhances our confidence level to decide whether a program is effective or not.

A summative evaluation of the a reading intervention program employed this multiple-matched sample procedure to investigate whether or not students enrolled in this program gained any achievement advantage over students who were not enrolled in this program (Yang, Li, Modarresi & Tompkins, 2004). Five hundred replicated matched samples were used in this evaluation. Figure 2, shown below, was the distribution of those effect size measures, across 500 replicated comparisons. The negative effect results indicated in this distribution seemed to suggest that this reading intervention program did not raise reading performance of students who were enrolled in this program, compared with similar groups of students.

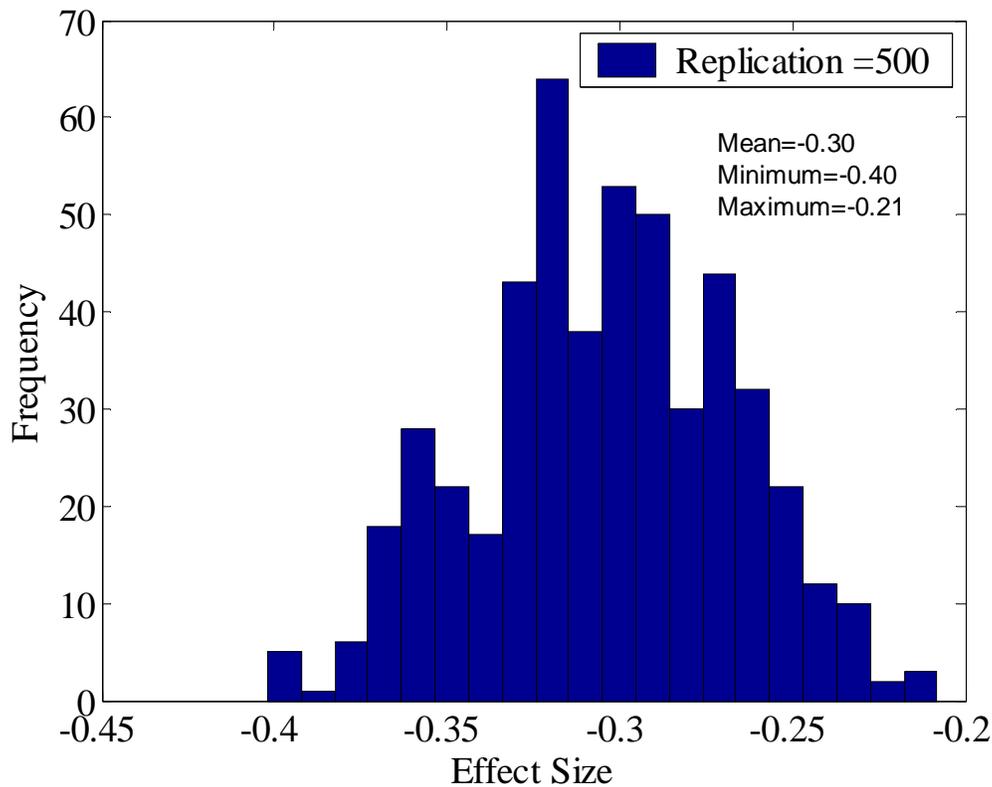


Figure 2. The Frequency of Effect Sizes for the Reading-Intervention Group

### C. The Choice between the Multiple and Single Matched-Sample Approach

After the approach of crating multiple sets of matched samples was developed, another issue that has never been addressed is: Is the results produced by the single-matched-sample approach reliable enough to make a right decision (e.g., phase out or keep a program) for the policy-decision makers, compared to the multiple-matched-sample approach?

Based on the results from our empirical studies, the results produced by the single-matched-sample approach are not reliable enough to make a right decision for the case with a small sample size (e.g., 30, 50). For example (refer the Figure 3 below) for the Music & Technology program, the effect size measure was about 0.27 based on the study (Yang, Li, Modarresi and Tompkins, 2003). The results of that study was yielded by a single matched-sample approach. This program was not phased out by the policy-decision makers partly because the effect size value was larger than 0.2.

When the 200 multiple-matched approach was used again, the mean of 200 effect size measures was  $-0.16$ . This latter finding seems to contradict with the previous finding.

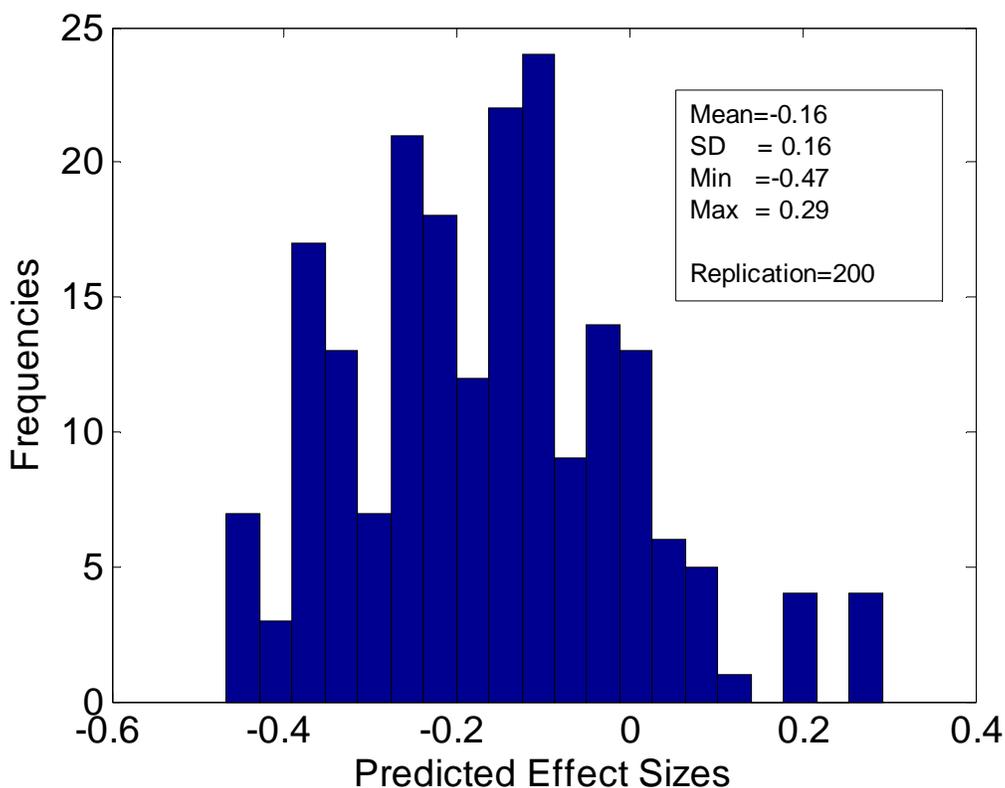


Figure 3. The Frequency of Effect Sizes for the Music and Technology Program Group

## VI. Discussions and Conclusions

### A. The Solution of the Zero-One Linear Programming

In reality, the group matching procedure may be done using linear programming software packages (e.g., LINGO, see <http://www.lindo.com>) by “feeding” the  $\mathbf{A}$  matrix and  $\mathbf{b}$  vector into those computer software programs. However, if the problem is too complicated as usually occurred in the matched procedure, it may not be very feasible to do so. For the individual matching procedure as well as the multiple-sets-of-matched-sample procedure, it almost cannot be done by directly using the software packages. Users need to write computer codes (e.g., C++ language, or MATLAB, The MathWorks, Inc., 2003) to call the callable libraries (e.g., LINDO API, LINDO Systems, Inc. 2003) to do so. For the solution presented in Table 3, the LINDO API was called into the MATLAB to seek the solution of the vector of  $\mathbf{x}$  in Equation 4. All solutions met all the constraints without any difficulties, as well as in a timely manner (e.g. less than a second per matching).

## **B. The Features of the *Zero-One Linear Programming Matching Method***

The *non-randomized comparison group pretest-posttest design* (refer to Shadish, Cook & Campbell [2002], p. 136) is one of the most appropriate and realistic summative evaluation designs in assessing the effectiveness efficacy of any programs. As indicated above, since the propensity-score-based matched procedure could not be sound for this evaluation design, another matched-sample procedure that utilizes the *Zero-One Linear Programming* approach was introduced in this paper. This approach could be more appropriate to this evaluation design due to its appealing features that are enunciated below:

Compared to the existing propensity score matching method, this matching method does not require the choice of a statistical model, that are often used for the computation of propensity score. This prevents any negative consequence that might occur when any assumptions made for the selected model are seriously violated.

Moreover, this matching procedure can handle the covariate of the pretest score more appropriately and is very efficient in matching as many demographic and initial ability (or pretests) variables as the researcher desires. The identical distribution of different types of students (e.g., Male/White/Poverty, Female/White/Poverty, Male/Asian/Poverty, etc.) between the experimental and matched samples is a promising feature that can hardly be found in pre-existing matching procedures in the literature. One evaluation study shows 180 types of students to be perfectly matched between the experimental and the matched samples. Of course, this matching procedure can easily handle the case of more than 180 types of students to be matched ,if the required data are available.

Finally, creating multiple sets of matched samples that are then treated as replicated-similar multiple control groups is another appealing feature that other existing matching methods have never addressed.

## **C. Statistical Modeling Followed by the Matching Method**

### Evaluation for only One Program

After the matching procedure, a small pretest score difference between the treatment group and its matched sample remained. The ANCOVA can be used to control for the effects of the small pretest score difference. When the matching control is integrated with the ANCOVA analysis, ANCOVA resulted in *adjusted posttest means* for both groups under the constraint of two groups' pretest means being equal, as well as two groups' matching variables and the combinations of matching variables being equal. The latter constraint makes the ANCOVA-based adjusted means more defensible. If the data structure is hierarchical, the HLM model is preferred.

When the measurement error on the pretest score is taken into account, multiple matched samples can be generated and the performance of the treatment group can be compared with each of the multiple matched samples using the ANOVA analysis. The distribution (e.g., mean, minimum, and maximum values) of the effect size measure, across multiple replicated comparisons, can then be computed and used to assess the efficacy of

any program. This enhances our confidence level to decide whether a program is effective or not.

#### Evaluation for Multiple Programs at a Time

When multiple programs or schools (e.g., 30 schools) are evaluated simultaneously, a set of matched samples will be generated to serve as a control group for a specific program. Under the logic of the randomized block design (Kirk, 1995), since students' outcome scores are more likely to be homogeneous within each program than across programs, the data from each program and its matched sample can be treated as a block. Within each block are students that received a program treatment or those that received no program treatment. The data from all blocks can then be aggregated to form a factor called "block" in the statistical context. Randomized Block Design assumes that unit (or student) assignment into both the treatment and non-treatment groups is random within each block. Since student assignment was not random in this study, however, the use of matched samples for the non-treatment groups represented an attempt to correct for this problem. This quasi-randomized procedure will make the adjusted-mean difference between two groups (the experimental group and control group) interpretable. The combined use of ANCOVA and Quasi-Randomized-Block Design was designed to reduce the error variance so that a more precise estimate of a treatment effect could be obtained. However, if the pretest-test means differ widely in different programs, using the ANCOVA to analyze each program's data separately, instead of using this combined method to analyze all program's data simultaneously, is preferred because a separate ANCOVA analysis for each pair of treatment-and-matched data can avoid the use of the extreme extrapolation.

This combined method was used to investigate the effects of magnet school programs on the reading and mathematics performance of students in a school system (Yang, Li, Modarresi and Tompkins, 2003), in which the data collected from eight magnet programs were simultaneously analyzed and the effect sizes were then simultaneously computed for eight programs. Of course, HLM is another preferred method to model this type of data if the data structure meets the HLM requirements. The multiple-set-of matched-sample approach can also be applied here.

#### **D. Concerns of Matched Samples**

The relative success in creating a matched sample relies on which variables to base the matching as well as the selection of pretest scores to be minimized. In general, when more demographic variables are used in the process of matching, the result of the matched group's background is more similar to the experimental group; however, the gap of pretest score difference between the two groups might increase as the use of matching variables increases. This issue might be resolved by trying different combinations of matching variables and to see how large differences in the pretest score vary under different conditions. Based on those trial results, researchers then choose one suitable solution that fits their research interest the best.

In addition, the degree of successfully creating a matched sample also relies on whether the distributions of both the quasi-experimental group and its respective non-treatment group on the matching variables (especially on the pretest scores) substantially overlap or not. If both groups have more overlapping distributions on those matching variables, then the matched sample can be adequately obtained without the need of

selecting members from extreme tails of the distributions. For example, the Non-Magnet population might have more overlapping distributions if such a population is composed of more members who are eligible for a specific magnet program, but they are not placed in this magnet program due to some circumstances (e.g., schedule conflict, no intention to attend, etc.). In contrast, the Non-Magnet population might have less overlapping distributions if such a population is only composed of members who are not eligible at all for this magnet program. When the later scenario occurs, examination of the overlap of the two distributions will help alert researchers to the possibility of the regression effect among the matches (Shadish, Cook, & Campbell, 2002, p 121).

Finally, Shadish et. al. (2002) pointed out that matching can be done only on observed measures, so hidden bias may remain. Researchers should always be aware that in drawing conclusions from quasi-experimental designs, incorporated with the matched-sample method, causality may not be inferred due to the lack of random assignment of students to the treatment and matched groups. Only through random assignment of subjects can the two groups of subjects be equal on all possible observed and “hidden” variables. Of course, without a large sample, bias may remain even though random assignment is fully implemented.

## References

- Bryk, A. S., & Raudenbush, S. W. (1992). *Hierarchical linear models: Applications and Data Analysis Methods*. Newbury Park, CA: Sage Publications, Inc.
- CTBS/McGraw-Hill (1997). *Teacher's guide to TerraNova*. Monterey, CA: McGraw-Hill Companies, Inc.
- Kirk, R.E. (1995). *Experimental design: Procedures for the behavioral sciences*. Brooks/Cole Publishing Company, New York.
- Isaac, S. & Michael, W. (1995). *Handbook in research and evaluation*, (3<sup>rd</sup> Ed.). EdITS / Educational and Industrial Testing Service, C.A.
- Judd, C. M., Smith, E. R., & Kidder, L. H. (1991). *Research Methods in Social Relations*. San Francisco: Holt, Rinehart, and Winston, Inc.
- Li, Y. H. & Schafer, W. D. (2005a). Increasing the homogeneity of CAT's Item-exposure rates by minimizing or maximizing varied target functions while assembling shadow tests, *Journal of Educational Measurement*.
- Li, Y. H. & Schafer, W. D. (2005b). Trait parameter recovery using multidimensional computerized adaptive testing in reading and mathematics. *Applied Psychological Measurement*, 29, 1-23.
- LINDO Systems, Inc., (2003). *LINDO API: User's Manual*. LINDO Systems, Inc, Chicago, IL.
- Modarresi, S., Yang, Y. N. & Bulgakov-Cooke, D. & Li (2004, November). An investigation of the effects of an Algebra intervention program, *PLATO*, on the Algebra performance of students. Paper presented at the annual meeting of American Evaluation of Association, Atlanta, GA, November, 2004.
- Rosenbaum, P. R., & Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70, 41-45.
- Rosenbaum, P. R., & Rubin, D. B. (1984). Reducing bias in observational studies using subclassification on the propensity score. *Journal of the American Statistical Association*, 79, 561-524.
- Rosenbaum, P. R., & Rubin, D. B. (1985). Constructing a control group using multivariate matched sampling that incorporate the propensity score. *The American Statistician*, 39, 33-38.
- Rosenbaum, P. R. (1995). Quantiles in nonrandom samples and observational studies. *Journal of the American Statistical Association*, 90, 1424-1431.
- Rubin, D. B., Stuart, E. A., & Zanutto, E. L. (2004). A potential outcomes view of value-added assessment in education. *Journal of Educational and Behavioral Statistics*, 29(1), 103-116.
- Shadish, W. R., Cook, T. D., & Campbell, D. T. (2002). *Experimental and Quasi-experimental Designs for Generalized Causal Inference*, Houghton Mifflin Company, M.A.: Boston.
- The MathWorks, Inc. (2003). *MATLAB (Version 6.5): The language of technical computing [Computer program]*. Natick MA: The MathWorks, Inc.
- Theunissen, T. J. J. M. (1985). Binary programming and test design, *Psychometrika*, 50, 411-420.
- Theunissen, T. J. J. M. (1986). Optimization algorithms in test design, *Applied Psychological Measurement*, 10, 381-389.
- Thompson, B. (2002). "Statistical," "Practical," and "Clinical": How many kinds of significance do counselors need to consider?

- van der Linden, W. J., & Boekkooi-Timminga, E. (1989). A maximin model for test design with practical constraints, *Psychometrika*, 54, 237-247.
- Yang, Y.N., Li, Y. H., Modarresi, S., & Tompkins, L., J. (2003). An investigation of the effects of magnet school programs on the Reading and Mathematics performance of students. Prince George's County Public Schools, Maryland.
- Yang, Y.N., Li, Y. H., Modarresi, S., & Tompkins, L., J. (2004). An investigation of the maintenance effect for first graders enrolled in the reading intervention program on their succeeding grade two Reading performance. Prince George's County Public Schools, Maryland.
- Yang, Y.N., Li, Y. H., Modarresi, S., & Tompkins, L., J. (2005, April). Using the multiple-matched-sample and statistical controls to examine the effects of magnet school programs on the reading and mathematics performance of students. Paper presented at the annual meeting of the American Educational Research Association, Montreal, Canada.