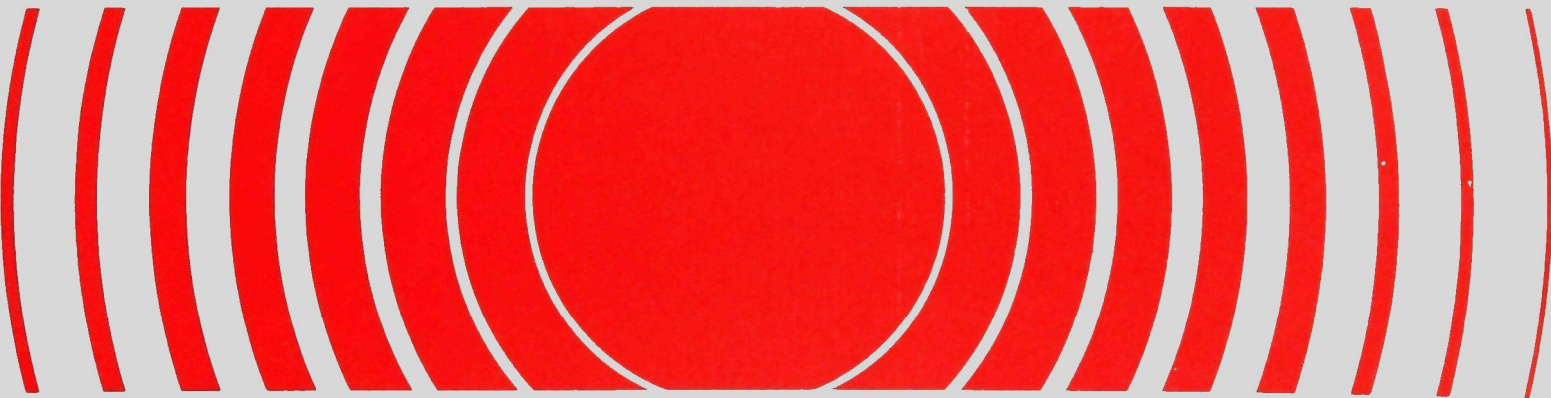Radiation

&EPA

# Confounding and Selection Bias in Case Control Studies

# Confounding and Selection Bias
# in
# Case Control Studies

Roderick J. A. Little

Paul R. Rosenbaum

January 1981

# Abstract

In case-control studies, the role of adjustments for bias, and in particular the role of matching, has been extensively debated. However, the absence of a formal statement of the problem has led to disagreements, confusion, and occasionally to erroneous conclusions. This paper formulates precisely and answers the following questions.

1) When is it necessary to adjust for a variable Z?

2) Given that the data analysis will adjust for the variable Z, is matching on Z the most efficient method of selecting controls?

In answering these questions, we draw a sharp distinction between bias caused by confounding in the population and bias caused by the method used to select the sample.

Acknowledgment

# CONTENTS

# 1.  Introduction

In case-control studies, the role of adjustments for bias, and in particular the role of matching, has been extensively debated (1-7). However the absence of a formal statement of the problem has led to disagreements, confusion, and occasionally to erroneous conclusions. In this paper we formulate precisely and answer the following questions:

a) When is it necessary to adjust for a variable Z ?

b) Given that the data analysis will adjust for the variable Z, is matching on Z the most efficient method of selecting controls?

## 2. Conditions Under Which Adjustment Is Necessary

### 2.1. Introduction

For simplicity we first consider measurement of the relationship between a disease (D) and an agent under study (E) in the presence of a single confounding variable Z. Extensions to the more realistic case where a set of variables are candidates for adjustment are outlined in Section 2.4.

We first seek a valid measure of association in the population of cases and controls, irrespective of the method of sampling. We then ask whether the sample estimate of this measure of association is a satisfactory estimate of the population quantity, that is, whether the sample estimate is not subject to selection bias.

## 2.2 Measures of Association in the Population

In a case-control study, the association between an agent (E) and disease (D) in the absence of confounding factors is measured by the population odds ratio

$$r = \frac{p(d|e)p(\overline{d}|\overline{e})}{p(\overline{d}|e)p(d|\overline{e})} = \frac{p(e|d)p(\overline{e}|\overline{d})}{p(\overline{e}|d)p(e|\overline{d})}, \tag{1}$$

where $D = d$ denotes disease, $D = \overline{d}$ denotes no disease, $E = e$ denotes exposure to the agent, $E = \overline{e}$ denotes no exposure to the agent, and $p(a|b)$ denotes the conditional probability that $A = a$ given $B = b$ in the population.

The relative risk

$$r^* = p(d|e)/p(d|\overline{e}) \tag{2}$$

is in some ways a more satisfactory measure of the effect of the agent. However the odds ratio approximates the relative risk if probability of disease is low, and unlike the relative risk it can be estimated from a case control study (8).

We now introduce a confounding factor Z, and suppose that a more appropriate measure of association is the adjusted odds ratio at $Z = z$,

$$r(z) = \frac{p(d|e,z)p(\overline{d}|\overline{e},z)}{p(\overline{d}|e,z)p(d|\overline{e},z)} = \frac{p(e|d,z)p(\overline{e}|\overline{d},z)}{p(\overline{e}|d,z)p(e|\overline{d},z)} \tag{3}$$

which approximates the relative risk at Z if the risk of disease is low in that subgroup of the population with $Z = z$. Note that in general $r(z)$ varies according to the value of z, and thus represents a set of measures of association.

If the population parameters r and r(z) are equal for all z, i.e.,

(*)  $r = r(z)$  for all z,

then the population relationship between disease (D) and exposure (E) is not confounded by z; otherwise the population relationship is confounded. The theorem below gives an expression for r in terms of r(z), and the subsequent discussion gives conditions under which confounding is absent, i.e. under which (*) holds. If confounding is present, then r and r(z) may yield strikingly different impressions concerning the effect of exposure on disease, and in this case, the choice of parameter to be estimated must depend on either assumptions or outside evidence concerning the biological mechanism that causes the disease.

Theorem

The adjusted and unadjusted odds ratios are related by the expression

$$r = (1 + b(z))r(z)$$

where

$$b(z) = \frac{p(z|\overline{d},e)p(z|d,\overline{e})}{p(z|d,e)p(z|\overline{d},\overline{e})} - 1 \qquad (4)$$

(Note: If z is continuous rather than discrete, p(z|d,e) is the probability density function of Z given D = d, E = e.)

Proof of Theorem

By Bayes' Theorem,

$$p(a|b,z) = \frac{p(a|b)p(z|a,b)}{p(z|b)},$$

and applying this expression in the formula (3) for r(z) leads to equation (3). //

In view of equation (4) we define b(z) to be the <u>relative confounding</u> <u>bias</u> of r at Z = z.  For example, if b(z) = 0.1 then the unadjusted odds ratio r deviates from the adjusted odds ratio at z by ten percent.  Two situations where the confounding bias is zero are of particular interest.  By inspection of equation (4), b(z) = 0 if either

(C1)      Z and D are conditionally independent given E, or

(C2)      Z and E are conditionally independent given D.

In the case where Z is categorical, these conditions are a special case of the well known collapsing theorem for contingency rables.  (See, for example, Bishop, Fienberg and Holland, (9), Section 2.4)

Conditions (C1) and (C2) are not the same as the condition proposed by Miettinen (3) under which adjustment is unnecessary, namely

(C2')  Z and E are independent.

The following example illustrates the difference between C2 and C2'.


<u>Example 1.</u>  The U.S. Environmental Protection Agency received a proposal to study the relationship between lung cancer (D) and radon$^{222}$ in well water (E).  Radon gas is released into the air when radon bearing well water is used in the home, for example, in showering.  There is some concern that as homes are made energy efficient and the rate of air exchange decreases, the concentration of radon daughters may increase in homes supplied with radon bearing well water.  The proposal contained a plan for a pilot study to determine whether well water radon levels (E) are independent  of smoking history (Z), an important confounding variable; i.e. to determine whether C2' holds.  If radon and smoking appear independent as a result of the survey, then the proposal would ignore smoking history.

However, our theorem shows that the relevant condition in deciding whether to adjust for smoking is not independence of smoking and radon but independence of smoking and radon within the diseased and non-diseased groups. Table 1 shows a (strictly hypothetical) population where radon level and smoking are marginally independent, so condition (C2') holds, but confounding is present because the adjusted odds ratios for radon and cancer are radically different in the smoking and non-smoking groups. The unadjusted odds ratio lies between these values, but is a poor summary of the relationship between radon and cancer for this population. Table 2 gives another hypothetical population where radon and smoking are unrelated within diseased and non-diseased groups (condition C2), so confounding is absent, but condition C2' does not hold.

Table 1.  Distribution of Radon (E), Lung Cancer (D) and Smoking (Z) in a

Hypothetical Population with a) E and Z independent and b) Unequal Odds Ratios.


Lung Cancer (D)

|  | D | | | | D | | | | D | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  | d | $\overline{d}$ | total | | d | $\overline{d}$ | total | | d | $\overline{d}$ | total |
| Radon (E) e | 20 | 5,980 | 6,000 | e | 180 | 3,820 | 4,000 | e | 200 | 9,800 | 10,000 |
| $\overline{e}$ | 110 | 23,890 | 24,000 | $\overline{e}$ | 190 | 15,810 | 16,000 | $\overline{e}$ | 300 | 39,700 | 40,000 |

odds ratio = .73        odds ratio = 3.92        odds ratio = 2.70

$Z = \overline{z}$:nonsmokers     $Z = z$:smokers        $Z = z$ or $\overline{z}$
                                                        smokers and nonsmokers


Table 2.  Distribution of Radon (E), Lung Cancer (D) and Smoking (Z)

in a Hypothetical Population with a) E and Z independent given D,

and hence b) Equal Odds Ratios


Lung Cancer (D)

|  | D | | | | D | | | | D | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  | d | $\overline{d}$ | total | | d | $\overline{d}$ | total | | d | $\overline{d}$ | total |
| e | 40 | 8,000 | 8,040 | e | 160 | 2,000 | 2,160 | e | 200 | 10,000 | 10,200 |
| Radon (E) $\overline{e}$ | 60 | 32,000 | 32,060 | $\overline{e}$ | 240 | 8,000 | 8,240 | $\overline{e}$ | 300 | 40,000 | 40,300 |

odds ratio = 2.67        odds ratio = 2.67        odds ratio = 2.67

$Z = \overline{z}$            $Z = z$              $Z = z$ or $\overline{z}$

Example 2.

We have seen that either condition (Cl) or (C2) implies that the confounding bias is zero. If Z is binary, it is easily shown that the converse holds, that is, a confounding bias of zero implies either (Cl) or (C2). However if Z has more than two categories, then populations can be constructed where neither (Cl) nor (C2) are satisfied and yet the confounding bias is still zero. An example for trichotomous Z is given in Table 3. It is readily verified that the adjusted odds ratios all equal the unadjusted odds ratio (to within some rounding error), even though each pair of variables is neither conditionally nor marginally independent. Such examples are curiosities, and the two independence conditions (Cl) and (C2) are more useful than equation (4) in practice.

Table 3. Hypothetical Population where a) Adjusted and Unadjusted Odds Ratios

of D and E Are Equal, b) D and Z Are Not Independent Given E,

and c) E and Z Are Not Independent Given D

a) D and E given Z

| | | D | | | | D | | | | D | | | | D | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | d | $\overline{d}$ | | | d | $\overline{d}$ | | | d | $\overline{d}$ | | | d | $\overline{d}$ |
| E | e | 30 | 20 | E | e | 150 | 59 | E | e | 120 | 21 | E | e | 300 | 100 |
| | $\overline{e}$ | 20 | 120 | | $\overline{e}$ | 28 | 99 | | $\overline{e}$ | 52 | 81 | | $\overline{e}$ | 100 | 300 |

Z = 1                Z = 2                Z = 3                Z = 1, 2 or 3

odds ratio = 9.0    odds ratio = 9.0    odds ratio = 8.9    odds ratio = 9.0

b) D and Z given E

| | | Z | | | | | | Z | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | | | | 1 | 2 | 3 |
| D | d | 30 | 150 | 120 | | D | d | 20 | 28 | 52 |
| | $\overline{d}$ | 20 | 59 | 21 | | | $\overline{d}$ | 120 | 99 | 81 |

E = e                        E = $\overline{e}$

c) E and Z given D

| | | Z | | | | | | Z | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | | | | 1 | 2 | 3 |
| E | e | 30 | 150 | 120 | | E | e | 20 | 59 | 21 |
| | $\overline{e}$ | 20 | 28 | 52 | | | $\overline{e}$ | 120 | 99 | 81 |

D = d                        D = $\overline{d}$

## 2.3. The Effects of Sample Selection

We have established conditions under which the unadjusted and adjusted odds ratios in the population are equal and therefore confounding is absent in the population. However these conditions are not sufficient for adjustment of the sample odds ratio to be unnecessary. The method of selection of cases and controls may be such that the unadjusted odds ratio for the sample is a biased estimate of its population analog. Adjustment may be necessary to eliminate (or at least to reduce) this bias.

To clarify conditions under which selection bias arises, it is convenient to introduce a sample indicator variable S, defined for each individual of the population, which takes value one if an individual is selected into the study and zero otherwise. The method of sampling can be characterized in terms of assumptions about the probability distribution of S given Z, D and E (Cf Rubin, 10). The following conditions are of particular interest since they characterize common methods of data collection:

(C3)  S is independent of Z, D and E.

(C4)  S is independent of D and E, given Z.

(C5)  S is independent of D and Z, given E.

(C6)  S is independent of D, given Z and E.

(C7)  S is independent of E and Z, given D.

(C8)  S is independent of E, given Z and D.

Conditions (C3) and (C4) correspond to randomized experiments where individuals are selected at random from the population and values of D and E are measured. Conditions (C5) and (C6) underlie cohort studies

if individuals are selected at random within exposed and non-exposed groups, and values of D are measured. Conditions (C7) and (C8) underlie case control studies where individuals are selected at random within diseased and non-diseased groups, and values of E are recorded. The odd numbered conditions (C3, C5, C7) correspond to situations where Z is not used as a stratifying variable for data collection; in particular, matching on Z has not taken place. The variable Z is recorded for the analysis. The even numbered conditions (C4, C6, C8) correspond to situations where Z is used as a stratifying variable, for example, by matching cases and controls on Z.

A key aspect of these conditions is that they imply random sampling within the indicated groups. In observational studies this assumption is subject to doubt since the sampling of cases and/or controls is not entirely controlled by the researcher. We shall return to this point later.

Since the sample adjusted odds ratio $\hat{r}_s(Z)$ is calculated from the selected individuals, all of whom have $S = 1$, it estimates the population adjusted odds ratio conditional on $S = 1$, that is,

$$r_s(z) = \frac{p(d \mid e,z,s=1)p(\bar{d} \mid \bar{e},z,s=1)}{p(\bar{d} \mid e,z,s=1)p(d \mid \bar{e},z,s=1)} = \frac{p(e \mid d,z,s=1)p(\bar{e} \mid \bar{d},z,s=1)}{p(\bar{e} \mid d,z,s=1)p(e \mid \bar{d},z,s=1)}$$

Hence the sample adjusted odds ratio estimates the population adjusted odds ratio if and only if $r_s(z) = r(z)$ for all z. Applying the argument in the proof of the theorem, we can write

$$r_s(z) = r(z) \ (1+b_s(z)),$$

where

$$b_s(z) = \frac{p(s=1 \mid d,e,z)p(s=1 \mid \overline{d},\overline{e},z)}{p(s=1 \mid \overline{d},e,z)p(s=1 \mid d,\overline{e},z)} - 1.$$

Accordingly we define $b_s(Z)$ to be the <u>relative selection bias</u>* of the sample adjusted odds ratio, $\hat{r}(Z)$. The relative selection bias is zero if any of the conditions (C3) to (C8) for the selection process is satisfied. Hence the sample adjusted odds ratio is not biased for clinical trials, prospective or case/control studies, provided the appropriate random sampling condition (C3), ..., or (C8) can be justified.

Stronger conditions are required for the unadjusted sample odds ratio $\hat{r}_s$ to be free of selection bias. Let us suppose that the confounding is absent in the population so that $r$ is an appropriate measure of association between disease and exposure. The sample odds ratio $\hat{r}_s$ estimates the unadjusted odds ratio conditional on $S = 1$, that is,

$$r_s = \frac{p(d \mid e,s=1)p(\overline{d} \mid \overline{e},s=1)}{p(\overline{d} \mid e,s=1)p(d \mid \overline{e},s=1)} = \frac{p(e \mid d,s=1)p(\overline{e} \mid \overline{d},s=1)}{p(\overline{e} \mid d,s=1)p(e \mid \overline{d},s=1)}.$$

This parameter is related to $r$ by the expression

$$r_s = r(1 + b_s),$$

where

$$b_s = \frac{p(s=1 \mid d,e)p(s=1 \mid \overline{d},\overline{e})}{p(s=1 \mid \overline{d},e)p(s=1 \mid d,\overline{e})} - 1.$$

Hence we define $b_s$ to be the relative selection bias of $\hat{r}_s$. It is zero if any one of the conditions (C3), (C5) or (C7) are satisfied, but is not

---

* Note that the selection bias has a slightly different form than the confounding bias, in that the values of D and E in the numerator and denominator have been switched.

in general zero if Z is controlled at the design stage of the study, that is, when conditions (C4, (C6) or (C8) apply. Hence, for example, matching at the design stage generally leads to a requirement to adjust at the analysis stage, even when the confounding bias is zero. Of greater importance is the fact that even when Z is not controlled in the selection process, there may still be a need for adjustment in the analysis, because the sample adjusted odds ratio estimates the population adjusted odds ratio under weaker conditions (e.g. C8) on the selection process than are required for the sample unadjusted odds ratio to estimate the population unadjusted odds ratio.

## 2.4 More than One Covariate.

In practice, a number of confounding factors are usually present in the design and analysis of a study, and thus a more realistic problem is whether to adjust for a covariate Z in addition to a set of other confounding variables $U = (U_1, \ldots, U_k)$. The previous arguments are easily extended to this case by conditioning throughout on variables U. The odds ratio $r(Z)$ adjusted for Z is replaced by the odds ratio $r(Z,U)$ adjusted for Z and U. The sample version of the adjusted odds ratio estimates

$$r_s(z,u) = r(z,u)(1 + b_s(z,u))$$

with relative selection bias

$$b_s(z,u) = \frac{p(s|z,u,d,e)p(s|z,u,\overline{d},\overline{e})}{p(s|z,u,\overline{d},e)p(s|z,u,d,\overline{e})} - 1.$$

In particular, this bias is zero when S is independent of D given Z,U,E or S is independent of E given Z,U,D. The population odds ratio $r(u)$ adjusted for U is $r(u) = r(u,z)(1 + b(z|u))$ with relative confounding bias

$$b(z|u) = \frac{p(z|\overline{d},e,u)p(z|d,\overline{e},u)}{p(\overline{z}|d,e,u)p(z|\overline{d},\overline{e},u)} - 1,$$

the bias being zero when Z is independent of D given E,U, or when Z is independent of E given D,U. The sample odds ratio $r_s(u)$ adjusted for u is

$$r_s(u) = r(u)(1 + b_s(u))$$

with relative selection bias

$$b_s(u) = \frac{p(s|u,d,e)p(s|u,\overline{d},\overline{e})}{p(s|u,\overline{d},e)p(s|u,d,\overline{e})} - 1,$$

and in particular the bias is zero when S is independant of D given U,E, or when S is independent of E given U, D.  The counter example to Miettinen's conditions described by Fisher and Patil (6) fails to satisfy the condition that the relative confounding bias is zero, which explains why adjustment is necessary in their case.

## 3. Does Matching Increase Power?

Now we ask: Given that the analysis will adjust for a variable Z, does matching on Z in the design increase power? That is, does matching on Z increase the probability of detecting a real association between disease D and exposure E, adjusting for Z?

We suppose the variable is categorized with I levels, and thus divides the population into I strata. There are $N_i$ ($i=1,\ldots,I$) cases available in the $i^{th}$ stratum. We plan to use all $N = \sum_{i=1}^{I} N_i$ available cases in the case-control study, and to select a total of M controls for comparison. The question is how to best choose the number $M_i$ of controls in the $i^{th}$ stratum, subject to the condition $\sum_{i=1}^{I} M_i = M$. Thus $N_i$, N and M are fixed; the $M_i$'s are to be chosen.

By definition, frequency matching of cases and controls takes

$$M_i = kN_i$$

with

$$k = \frac{M}{N} \quad.$$

Let

$P_{1i}$ = population proportion of cases exposed in stratum i.

$P_{2i}$ = population proportion of controls exposed in stratum i.

$\delta_i = P_{1i} - P_{2i}$

$P_i = (P_{1i} + P_{2i})/2$

and let $\hat{P}_{1i}$, $\hat{P}_{2i}$, $\hat{\delta}_i$ and $\hat{P}_i$ denote the corresponding sample quantities.

The null hypothesis $H_0: P_{1i} = P_{2i}$ for $i = 1,\ldots,I$ is equivalent to the null hypothesis that the adjusted odds ratio of D and E given Z is zero for all values of Z.

The statistic

$$C = \frac{\displaystyle\sum_{i=1}^{I} \frac{N_i M_i}{N_i + M_i} \hat{\delta}_i}{\sqrt{\displaystyle\sum_{i=1}^{I} \frac{N_i M_i}{N_i + M_i} \hat{P}_i(1-\hat{P}_i)}}$$

may be used to test this hypothesis. In moderate to large samples, the test based on C is nearly equivalent to those of Cochran (11), Mantel-Haenszel (12) and Birch (13), but is easier to manipulate in the current problem.

The asymptotic expectation of C is

$$E_A(C) = \frac{\displaystyle\sum_{i=1}^{I} \frac{N_i M_i}{N_i + M_i} \delta_i}{\sqrt{\displaystyle\sum_{i=1}^{I} \frac{N_i M_i}{N_i + M_i} P_i(1-P_i)}} \tag{1}$$

We find $M_i$ to maximize (1) subject to the constraint $M = \sum_i M_i$. Since the nonull variance of C is nearly 1, and since C is asymptotically normal, maximizing $E_A(C)$ is nearly equivalent to maximizing the asymptotic power.

Differentiating the log of (1) subject to the constraint $\sum_i M_i = M$ yields

$$\frac{d \log_e C}{dM_i}$$

$$= \frac{d}{dM_i} \left[ \log_e \sum \frac{N_i M_i}{N_i M_i} \delta_i - 1/2 \log_e \sum \frac{N_i M_i}{N_i M_i} P_i(1-P_i) - \lambda \left( \sum M_i - M \right) \right]$$

$$= \left[ \frac{N_i}{N_i + M_i} \right]^2 \left[ \frac{\delta_i}{\sum \frac{N_i M_i}{N_i M_i} \delta_i} - \frac{P_i(1-P_i)}{2 \sum \frac{N_i M_i}{N_i M_i} P_i(1-P_i)} - \lambda \right]$$

Cochran (11) observed that if the odds ratio is constant over strata then $\delta_i / P_i(1-P_i)$ is nearly constant. Assuming $\delta_i / P_i(1-P_i)$ is constant, we find the optimal allocation $M_i$ satisfies

$$\frac{N_i}{N_i + M_i} = \sqrt{\frac{2\lambda \sum \frac{N_i M_i}{N_i M_i} \delta_i}{\delta_i}}$$

$$\propto \sqrt{\frac{1}{\delta_i}}$$

or equivalently

$$\frac{N_i}{N_i + M_i} = \sqrt{\frac{2\lambda \sum \frac{N_i M_i}{N_i M_i} P_i(1-P_i)}{P_i(1-P_i)}} \propto \sqrt{\frac{1}{P_i(1-P_i)}}$$

If $P_{1i} = P_1$, $P_{2i} = P_2$, for all $i$, then both $\delta_i = P_1 - P_2$ and $P_i(1-P_i)$ are constant, and frequency matching is optimal. Otherwise, still assuming the odds ratio is constant, the optimal allocation takes more controls ($M_i$ larger) from strata with a larger difference $\delta_i$ in exposure proportions, or equivalently, with a larger variance $P_i(1-P_i)$.

# References

1. Miettinen, O.S. The matched pairs design in the case of all-or-none responses. Biometrics, 1968, 24:339-352.

2. Bross, I.D.J. How case-for-case matching can improve design efficiency. Amer. J. Epid. , 1969, 89:359-363.

3. Miettinen, O.S. Matching and design efficiency in retrospective studies. Amer. J. Epid., 1970, 91:111-118.

4. Hardy, R.J., White, C. Matching in retrospective studies. Amer. J. Epid., 1971, 93:75-6.

5. Seigel, D.G., Greenhouse, S.W. Validity in estimating relative risk in case-control studies. J. Chron. Dis., 1973, 26:219-225.

6. Fisher, L. and Patil, K. Matching and unrelatedness. Amer. J. Epid., 1974, 100:347-349.

7. Miettinen, O.S. Confounding and effect modification. Amer. J. Epid., 1974, 100:350-353.

8.  Cornfield, J.  A method of estimating comparative rates from clinical data.  J. Natl. Cancer Inst., 1951, 11:1269-1275.

9.  Bishop, Y.M.M., Fienberg, S.E., Holland, P.W.  Discrete Multivariate Analysis.  Cambridge, Massachusetts:  MIT Press, 1975.

10.  Rubin, D.B.  Inference and missing data.  Biometrika, 1976, 63:581-592.

11.  Cochran, W.G.  Some methods for strengthening the common chi square tests.  Biometrics, 1954, 10:417-451.

12.  Mantel, N., Haenszel, W.  Statistical aspects of the analysis of data from retrospective studies of disease.  J. Natl. Cancer Inst., 1959, 22:719-748.

13.  Birch, M.W.  The detection of partial association, I:  the 2x2 case.  J. Royal Statistical Society, 1964, series B, 26:313-324.

# TECHNICAL REPORT DATA
*(Please read Instructions on the reverse before completing)*

| 1. REPORT NO.<br>EPA 520/8-81-004 | 2. | 3. RECIPIENT'S ACCESSION NO. |
|---|---|---|
| 4. TITLE AND SUBTITLE<br>Confounding and Selection Bias in Case Control Studies | | 5. REPORT DATE<br>January 1981 |
| | | 6. PERFORMING ORGANIZATION CODE |
| 7. AUTHOR(S)<br>Roderick J. A. Little<br>Paul R. Rosenbaum | | 8. PERFORMING ORGANIZATION REPORT NO. |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS<br>Office of Radiation Programs<br>U.S. Environmental Protection Agency<br>Washington, D.C. 20460 | | 10. PROGRAM ELEMENT NO. |
| | | 11. CONTRACT/GRANT NO. |
| 12. SPONSORING AGENCY NAME AND ADDRESS | | 13. TYPE OF REPORT AND PERIOD COVERED |
| | | 14. SPONSORING AGENCY CODE |

15. SUPPLEMENTARY NOTES

16. ABSTRACT

In case-control studies, the role of adjustments for bias, and in particular the role of matching, has been extensively debated. However, the absence of a formal statement of the problem has led to disagreements, confusion, and occasionally to erroneous conclusions. This paper formulates precisely and answers the following questions. 1) When is it necessary to adjust for a variable Z? 2) Given that the data analysis will adjust for the variable Z, is matching on Z the most efficient method of selecting controls? In answering these questions, we draw a sharp distinction between bias caused by confounding in the population and bias caused by the method used to select the sample.

17. KEY WORDS AND DOCUMENT ANALYSIS

| a. DESCRIPTORS | b. IDENTIFIERS/OPEN ENDED TERMS | c. COSATI Field/Group |
|---|---|---|
| biometry<br>epidemiologic methods<br>research design | | |

| 18. DISTRIBUTION STATEMENT<br><br>Unlimited | 19. SECURITY CLASS *(This Report)*<br>Unclassified | 21. NO. OF PAGES<br>27 |
|---|---|---|
| | 20. SECURITY CLASS *(This page)*<br>Unclassified | 22. PRICE |

EPA Form 2220-1 (Rev. 4-77)   PREVIOUS EDITION IS OBSOLETE