

**EVALUATING PROGRAMMES:
EXPERIMENTS,
NON-EXPERIMENTS
AND PROPENSITY SCORES**

**Denis Conniffe, Vanessa Gash
and
Philip J. O'Connell**

March 2000

Working Paper No. 126

Working Papers are not for publication
and should not be quoted without prior
permission from the author(s).

Evaluating Programmes: Experiments, Non-Experiments and Propensity Scores

DENIS CONNIFFE, VANESSA GASH and PHILIP O'CONNELL

The Economic and Social Research Institute, Dublin

Abstract: Evaluations of programmes - for example, labour market interventions such as employment schemes and training courses - usually involve comparison of the performance of a treatment group (recipients of the programme) with a control group (non-recipients) as regards some response (gaining employment, for example). But the ideal of randomisation of individuals to groups is rarely possible in the social sciences and there may be substantial differences between groups in the distributions of individual characteristics that can affect response. Past practice in economics has been to try to use multiple regression models to adjust away the differences in observed characteristics, while also testing for sample selection bias. The Propensity Score approach, which is widely applied in epidemiology and related fields, focuses on the idea that "matching" individuals in the groups should be compared. The appropriate matching measure is usually taken to be the prior probability of programme participation. This paper describes the key ideas of the Propensity Score method, compares it with the common approach in economics, reviews the arguments in the literature and illustrates application by reanalysis of some Irish data on training courses.

I INTRODUCTION

Application of the direct experimental approach in the economy and society is usually considered unpalatable, or even unethical, even when it would clearly provide the ideal comparison. For example, we would like to assess an active labour market policy - say, training to enhance skills - by drawing a large and fully random sample from the relevant population and then randomly assigning individuals to a training group and a control group. Then, although an individual's subsequent average performance (in terms of employment, earnings, productivity, or whatever) will depend on characteristics like age, education and previous work experience, these factors cancel out of the difference in the averages for the two groups¹. So the difference can be validly interpreted as the effect of the programme or policy.

While there have been a few such assessments in the US (for example, LaLonde, 1986), allocation to a control group can be seen as disadvantageous, so that randomisation is unpopular, to say the least.

Evaluations have sometimes been based on the performances of the programme participants only, without employing any control group. But as it is most unlikely that, in the absence of a programme, individuals would not have tried to improve their own positions, information has to be sought about this, or else the programme benefits could be considerably overestimated. An Irish example is

¹ The ideal comparison, though impossible to make, would use the *same* people to compare the effect of participation in the programme with non-participation. In experimental approaches the *causal* effect of a treatment on an individual is defined as the difference in the potential responses when receiving and not receiving the treatment. The average of these differences over the whole population is the parameter of interest. While theoretically a useful concept, it is usually unmeasurable and has to be assumed estimable by the difference between the means of the treatment and control groups, given prior randomisation of individuals to groups.

provided by Breen and Halpin (1988), who evaluated the FAS *Enterprise* programme by interviewing a sample of participants and, besides ascertaining how well they had got on, also asked what they would have done had the programme not existed. Again, Breen and Halpin (1989), in assessing a job subsidisation scheme, asked employers if they would have hired anyway in the absence of the scheme. In both these studies and, no doubt, in many others, there was simply no other way to proceed. But depending on questions of this nature, with all the possibilities for “wisdom by hindsight”, seems less attractive than comparing with a ‘control’ group, even if the allocations of individuals to groups has been far from random.

Situations when we have observational data on a programme (henceforward called treatment) group and on a control group, but without the deliberate randomisation of individuals to groups that characterises true experimentation, are often called “natural experiments” in the economics literature. Without randomisation, there may well be substantial differences between groups in the distributions of individual characteristics that affect performance (henceforward called response). Sometimes quite simple methods are used to analyse the data. For example, if responses of all individuals are measured at two time points, corresponding to before and after treatment for the treatment group, the “difference of differences” method may be employed. This just compares the mean improvement in response for the treatment group with the mean improvement (if any) for the control group. The idea is that in taking “before” from “after” individual characteristics may cancel out. For example, this type of analysis was used by Eissa and Liebman (1996) in their study of labour supply response to an earned income tax credit.

However, most analyses of “natural experiments” have taken a far more sophisticated form. At simplest, a regression model, perhaps embodying economic theory and previous findings, is estimated relating the response to all the characteristics as well as to a dummy variable defining groups. This is interpreted as adjusting the treatment versus control comparison to what it would have been had there been no variation in characteristics between groups. The added possibility of sample selection biases is addressed by postulating another model expressing the probability that an individual is included, or not included, in the treatment group as a function of the individual’s characteristics. When estimated, this model *may* permit adequate adjustment to compensate for the biases. This econometric modelling

approach to the analysis of “natural experiments”², which will be reviewed in Section II, owes much to the work commenced by Heckman (1976, 1979) and continued in subsequent papers. However, for reasons to be discussed, the greater sophistication of approach does not always pay off in terms of elimination of bias, precision of estimation or clarity of findings .

The main purpose of this paper is to present an alternative approach based on propensity score matching. This has been developed in the statistical and biometrical, rather than in the econometric, literature. The propensity score for an individual is the probability that he or she belongs to the treatment group and is a function of the individual’s characteristics. ‘Matching’ is a simple device advocated by Cochran (1965, 1968) to help draw causative inferences from survey data. The Propensity Score method, which derives from the papers by Rosenberg and Rubin (1983a, 1984), will be detailed in Section III and contrasted with the modelling and Heckman adjustment approach. The approach already dominates in biomedical fields and many expository papers and reviews have appeared, including those by Drake and Fisher (1995), Rubin (1997), D’Agostino (1998) and Perkins, Tu, Underhill, Zhou and Murray (2000) as well as papers describing variants to methods (for example, Tu, Perkins, Zhou and Murray, 2000) and a host of papers describing applications to specific biomedical observational studies³.

There has been very little application of the propensity score method in the economic literature as yet, although there has been some debate in theoretical journals about methodological issues and whether variations on the approach, or even outright alternatives to it, might be more appropriate to the social sciences. We review this literature in Section IV and emphasise the degree of agreement that actually exists, in spite of apparently contentious matters. Finally, in Section V we illustrate the propensity score method by reanalysing data used by O’Connell and McGinnity (1997) to evaluate some training programmes.

II THE ECONOMETRIC, OR SPECIFICATION ERROR, APPROACH

If individuals from the population of interest were randomly allocated to treatment and control, we would be happy to accept the difference in mean responsiveness as a consistent estimate of the

² Some authors reserve the term “natural experiments” for cases where only simple analyses like the “difference of differences” are employed and call what is being described here the “modelling approach”. However, we think the defining feature is the presence of treatment and control groups.

³ However, the term “Natural Experiment” is not used outside the economic literature.

programme effect⁴. In a somewhat more roundabout account, we are fitting the model

$$y_{ij} = a + bD + e_{ij}, \quad (1)$$

where y_{ij} is the response (taken as continuous, for the present) of the j th individual in group i (with $i=1$ for the treatment group and $i=0$ for the control group), D is a dummy variable equal to 1 for the treatment group and to zero for the control group, a and b are intercept and coefficient respectively, and e_{ij} is the error or disturbance term. This term contains the effects of such characteristics as age, gender, education, etc. The true means of the control and treatment groups are a and $a+b$ respectively and so b is the true difference in means. Standard regression is appropriate because D and e_{ij} are independent (by randomisation), which is the fundamental condition for consistent estimation by Ordinary Least Squares. The OLS estimate of b is, as expected,

$$\bar{y}_1 - \bar{y}_0, \quad (2)$$

where the over bars denote means. The numbers in the groups do not, of course, need to be equal.

With non-randomised observational data D and e_{ij} in (1) are no longer independent and so (2) is no longer a plausible estimate. The approach in the econometric literature has then been to treat model (1) as “misspecified” because it does not contain all the covariates (of which there are p , say) that can affect the response. Instead

$$y_{ij} = a^* + b^*D + \sum_{k=1}^{k=p} c_k x_{kij} + u_{ij}, \quad (3)$$

where x_{kij} is the value of covariate k as measured on individual j of group i , is estimated. The asterisks on a and b just indicate they are somewhat different from what they were in equation (1), now that the covariates are included. Perhaps the OLS estimate of b^* is now an adequate measure of the treatment effect, having been “adjusted” for differences in covariate values, although there is still the possibility of self-selection bias. But before considering that topic, it is important to note the assumptions already being built into (3). The effects of covariates are being assumed linear and to operate identically in the two groups. The consequences of departures from assumptions on the

⁴ It should be said that some literature on programme evaluation questions the acceptability of this difference as the measure of programme effect and might even deny the appropriateness of the experimental paradigm. The points will be returned to, but are ignored for the present to simplify development.

estimate of b^* are far more serious when the distribution of covariate values differs greatly between groups than when it does not. Of course, it could be argued that sometimes economic theory or researchers' prior experience mean they "know" the true model. Or that the researcher may have tried out different functional forms, considered interactions of treatment with covariates etc. and, via goodness of fit measures and specification tests, found the "right" model to replace (3). But very often researchers are unsure about the correct model, a variety of equations may fit almost equally well (or poorly) and data may not be fully informative for various reasons, including multicollinearity. Even more seriously, some data points may be totally uninformative about the comparison of treatment and control and their inclusion just serves to magnify specification biases⁵.

Turning to self-selection bias, Heckman's (1979) formulation of the problem and its solution was also in terms of misspecification due to an omitted variable. Initially, Heckman just considered the drawing of one sample, with possible self-selection operational, from a population. So the situation was somewhat simpler than where two samples, corresponding to treatment and control, are involved. So, for the present, we rewrite equation (3) as if there is only one group, so that b , D and the subscript i do not appear. The model is now:

$$y_j = a + \sum_{k=1}^{k=p} c_k x_{kj} + u_j, \quad (4)$$

where u is assumed to have mean = 0 and variance = σ^2 . Given a truly random sample the coefficients could be estimated by OLS. Suppose the selection bias takes the form that sample y values are only observed if an unobserved, or latent, variable w is positive. However, x values are not only observed for the sample, but also for another sample (for which the y 's are not), or even for the whole population. (Many selection mechanisms can be reformulated into this framework – for example if y is a working mother's earnings and the x 's are age, years of education, etc., w could be the difference between earnings and perceived costs incurred by joining the workforce.) If w is independent of e there is still no difficulty, but it may well not be (in the example, the particularly high earners are more likely to exceed their perceived costs) and it may also depend on the x 's. Then u , instead of being random with mean zero, will have a component which is related to the x 's and the fundamental condition for OLS is broken. Heckman showed progress is possible if each w_j is assumed normally distributed with unit variance and mean

$$d_0 + \sum_{k=1}^{k=p} d_k x_{kj} = -z_j, \text{ say.} \quad (5)$$

Then, although w is only known to the extent of its sign, Heckman demonstrated that the true model is got by adding another variable λ_j , which is a function of z_j , to (4), giving

$$y_j = a^* + \sum_{k=1}^{k=p} c_k^* x_{kj} + g\lambda_j + u_j, \quad (6)$$

with
$$\lambda_j = \frac{\phi(z_j)}{\Phi(-z_j)}, \quad (7)$$

where ϕ is the ordinate (or density) and Φ is the integral (or distribution function) of the Standard Normal. Formula (7) is called the “inverse Mill’s ratio”, or sometimes “Heckman’s lambda”. As before, the asterisks in (6) just indicate that these coefficients differ from those in (4) in requiring interpretation as effects of covariates holding the new variable λ constant. Of course, the z ’s needed to calculate the λ ’s depend on the unknown d coefficients in (5) and this necessitates a two-stage estimation process. First a binary variable is defined, equalling 1 when y is observed and zero when it is not, and a probit analysis is conducted with the x ’s as explanatory variables to estimate the d ’s. Then the λ ’s are calculated and the coefficients of (6) estimated.

With this formulation of the selection bias problem, the observed covariates are being taken as quite sufficient for adequate model specification and estimation. The x ’s actually come into equation (6) twice – explicitly in the second term on the right hand side and implicitly in λ . It could be said that Heckman sees the effect of an omitted covariate w as replacing equation (4), which is linear in the x ’s, by an equation (6), which still involves only the x ’s, but which has the non-linear component λ . Indeed, estimation of (6) could be approached as a case of non-linear equation estimation. Note, however, the crucial importance of the assumption, embodied in (5), that the unobserved w has a linear regression on the observed x ’s and is normally distributed. These assumptions permit the x ’s to “handle” the omitted variable and identify the appropriate non-linear adjustment to the equation.

The situation when comparing treatment and control groups is just a little more complicated. Now individuals may select themselves into either the treatment or control group. For example, individuals of above average ability (the unobserved w could now be individual’s ability minus mean ability and is positive for the treatment group) might judge a training programme to be likely to benefit them and opt

⁵ This point will be returned to.

for training, while those of below average ability might judge themselves unable to benefit from training and opt for the control group⁶. The Heckman lambda values, λ_{ij} , for the treatment group are obtained from z_{ij} values derived from a probit analysis where the binary variable equals 1 for the treatment group and 0 for the control group and the x 's are the explanatory variables. They are given by (7). In the corresponding derivation for the control group, w is negative, the 1's and 0's of the probit analysis are interchanged, but the covariates are the same. So it is intuitively clear that the λ_{cj} values are given by

$$\lambda_{cj} = \frac{-\phi(-z_j)}{\Phi(z_j)} = \frac{-\phi(z_j)}{\Phi(z_j)}.$$

The modification of model (3) to be actually estimated is then

$$y_{ij} = a' + b'D + \sum_{k=1}^{k=p} c'_k x_{kij} + g\lambda_{ij} + u_{ij}, \quad (8)$$

where the primes indicate coefficients have changed from those of equation (3) because of inclusion of the lambda variable. A consistent estimate of b' will now follow from OLS⁷, assuming, of course, that all has been specified correctly in the model.

As in the sample from a population case, the approach is heavily dependent on the postulated selection bias process corresponding reasonably to reality and on the bivariate normality assumption about w and y . It is also sensitive to specification errors. For example, when the response equation should have contained some non-linear components, it is known that the λ_{ij} term in (8) can pick these up, suggesting there is selection bias where none exists. In other cases the λ_{ij} as estimated functions of the x 's can be nearly collinear with the covariates in the model, leading to unstable coefficients, high

⁶ Postulating such behaviour might seem to be attributing near prophetic powers to the individuals in both groups, since they have neither experienced the programme nor is any objective assessment of its effectiveness yet available. However, such postulated behaviour has also been used (for example in Heckman and Smith, 1996) to make radical criticisms of the standard experimental paradigm. Not only should the measure of treatment effect really include some selection effect, but randomisation could destroy an experiment, because high ability individuals would not accept assignment to a control group and by overt or covert means would attain the treatment group's conditions. However, Angrist, Imbens and Rubin (1996a) have incorporated "compliers" and "refusers" into the experimental framework.

⁷ Because some variance heterogeneity has been introduced at the probit analysis stage, OLS, while consistent, may be less efficient than the appropriate GLS or Maximum Likelihood solution. Many econometric computing packages now provide these procedures and we will not detail them here, but rather concentrate on the ideas and assumptions embodied in the econometric approach to the problem.

standard errors and insignificant “t” values. Then it can be difficult to conclude anything at all from the results after “adjusting” for selection bias. This sensitivity is reduced if some of the variables defining the mean of w_j in (5) do not occur (that is, are known to have zero coefficients) in the response equation (4). So some standard textbooks (for example, Johnston and Di Nardo, 1997, p. 450) recommend that Heckman correction should only be performed in these circumstances.⁸

It is worth noting that knowledge of zero coefficients (or “zero exclusions”, as they are termed) does permit a possible alternative analysis to the “standard Heckman” via (8). The zero exclusions define instrumental variables and so direct instrumental variable (IV) estimation is also possible. Returning to equation (3), the dummy variable D for treatment v control could now be regarded as endogenous. However, the IV method does not try, or need, to specify the precise source of endogeneity. Given instrumental variables, that is, some c 's zero, the endogeneity can be “washed out” by a procedure closely analogous to two stage least squares. No explicit assumptions about a latent variable or its distribution, or about the precise mechanism of the selection need be made. This may sometimes make the approach preferable to the “standard Heckman” analysis. Sample selection biases can easily be visualised as operating in a much more complex way than the scenario of people of high ability opting for training and of low ability opting for the control. There could be selection biases originating with the programme administrators, possibly interacting with, or countering, selection associated with individuals’ abilities. There could be selection effects at various stages of programmes – perhaps at recruitment to the programme before separation into treatment and control, then at the treatment allocation stage and perhaps selective dropout from the programme. Then it might be sensible to just think of an endogenous D as an aggregate effect of an assembly of selection biases and use the IV method.

The IV approach is still, of course, quite valid if the “standard Heckman” scenario holds. The high ability person opting for the training programme can be regarded as affecting both y and D , making the latter endogenous, which can be corrected by IV estimation. However, it could be argued that the IV method neglects information under these circumstances and is less efficient and less informative than the standard method. However, many authors have considered the assumptions required for Heckman correction of model (8) to be very fragile. Criticisms go back in the econometric literature as far as

⁸ If the response variable is itself binary - for example, if a programme is being evaluated in terms of employment gain – at least one zero coefficient is essential for any estimation.

Goldberger (1983), but have recurred strongly in recent propensity score and causal modelling literature, for example in Angrist et al (1996a) and Imbens and Rubin (1997). Empirical findings have sometimes differed with the approach to selection bias. The belief, following “standard Heckman” correction, that OLS estimation overestimates returns to education (because of selection on “ability”) has not been borne out by some IV analyses of US data (see Angrist and Kruger, 1991; Imbens and Rubin, 1997; and other references contained in the latter).

Using the term “standard Heckman” may actually be unfair to Heckman, whose views on the appropriate treatment of selection bias have evolved over time. Heckman and Robb (1986) and Heckman and MaCurdy (1986) were very positive about IV estimation and also examined other alternatives to the “standard Heckman” approach. Heckman and Hotz (1989), defending non-experimental studies of manpower training from criticisms about conflicting findings, emphasised that several selection bias adjustment procedures exist and that the correct choice depended on the source of the selection bias (see, especially, the reply to Holland, 1989). Heckman (1990) tried to relax the assumptions further⁹ and move to non- or semi-parametric approaches and some of his more recent papers will be considered in the next section. But it is the “standard Heckman” procedure that is in the textbooks and econometric software packages and that has been very frequently employed in programme evaluation and other assessments, such as on returns to education. Irish examples include Breen (1986); Breen (1991); Breen, Hannan and O’Leary (1995); O’Connell and Lyons (1995) and O’Connell and McGinnity (1997).

Returning to IV estimation, it is often very difficult to find truly valid and effective instrumental variables. Few economists can be fully confident in relying on economic theory to provide a set of “zero exclusions”. Inserting the candidate instrumental variables into the model, along with the treatment dummy variable and the “true” covariates, and finding their coefficients not statistically significant would provide some support. But the lack of statistical significance could be because the model did not fit well anyway, due to incorrect functional form, or important omitted covariates, or whatever. The choice and framing of assumptions justifying IV estimation, as presented in some econometric papers on programme evaluation and selection, have been subject to critical scrutiny by such authors as Little (1985), Angrist (1995) and Angrist et al (1996a). These criticisms have been disputed and there have even been disagreements over precisely what assumption is implied by a zero

⁹ So did Newey, Powell and Walker (1990)

exclusion. For example, Angrist et al (1996a) state that a valid instrumental variable z has to be *independent* of the response variable y given the non-zero coefficient variables (that is, independent of the disturbance term). Heckman (1996) maintained the required assumption is much weaker: z has only to be *uncorrelated* (mean independence) with the disturbance term. Probably not everyone would agree the assumption is much weaker, but in any event Angrist et al (1996b)¹⁰ replied that even if, without independence, there was no correlation with y as the response variable, there would be if $\log y$ (or another function) replaced y . Overall, recent work on IV stresses emphasis on ensuring the validity of instruments, although those that survive scrutiny are often found less related to the endogenous regressor of interest that might be desired (see, for example, Imbens and Rubin, 1997). These considerations will also be relevant in the next section when we discuss using IV in conjunction with propensity scores.

III THE PROPENSITY SCORE APPROACH

In the statistical analysis of designed experiments, there is rarely any attempt to estimate a “true” model for response in terms of *all* the variates that could affect it. Instead the model usually relates response to the factors being deliberately varied by the researcher and relies on a combination of seeking homogeneity of experimental material and randomisation to balance out all the other factors. For programme evaluation, this implies the treatment v control comparison is of paramount interest and that the covariate effects, if not subjected to randomisation, are just complicating nuisances. It might be argued that the relationship between response and some covariate is of interest in its own right. But that can probably be studied (and perhaps has) without the complication of conducting a comparative social experiment¹¹. Most causal findings in science have followed from controlled or randomised experiments and social scientists frequently complain they cannot conduct such trials. From this perspective, the problem with model (1) is not that covariates require specification, but the lack of randomisation to $D=1$ and $D=0$. In a sense, the *entire* problem is specification bias in allocation between treatment and control.

The propensity score solution depends on the idea of “matching” individuals from the treatment and control groups. Cochran (1968) gave the example of mortality rates for US smokers being lower, on

¹⁰ This argument appears elsewhere in the literature, too, for example in Imbens and Rubin (1997).

¹¹ Matters may be different if a particular treatment by covariate interaction (for example, what training achieves for females) is considered vital to the evaluation. But this is really a definition (which should have preceded the study) of a key sub-population and a stipulation that there be two treatment groups and two controls.

average, than for non-smokers – the reason being that smokers were younger, on average, than non-smokers. When groups of smokers and non-smokers of equal ages were compared, the mortality rates were always higher for smokers. Cochran advocated seeking causative effects from observational data by matching individuals from the treatment and control using all the covariates. If no important covariate has been omitted, it seems plausible to suppose that the difference between the responses of two such matching individuals, one receiving the treatment and the other the control, is the treatment effect plus a random element. Then averaging over the set of differences estimates the treatment effect. Cochran showed that perfect matching, in terms of exact equality of continuous covariates, is unnecessary and that matching on intervals can work well. Nonetheless the method will meet trouble if there are a lot of covariates, because the number of matching cells increases exponentially with the number of covariates and cells could quickly become empty of treatment individuals, or control cases, or both. That difficulty could be overcome, however, if all covariates could somehow be combined into a single “balancing score”. Several ways of constructing such a balancing score have been proposed, but Rosenberg and Rubin’s (1983a) “Propensity Score” approach is overwhelmingly the most popular.

The propensity score for an individual is the a priori probability (which is a function of the covariate values) that the individual is in the treatment group. Rosenberg and Rubin show that if we consider two sets of people, one set in the treatment group and one in the control group, with the same value of the propensity score, then the two sets have the same distributions of covariates. Rosenberg and Rubin gave a rigorous version of the intuitive argument that follows. If two individuals, one in the treatment group and one in the control group, have the same propensity score, their subsequent “allocation” to treatment or control can be regarded as random. It may as well have been decided by a coin toss. The difference in their responses is the treatment effect plus a random element and averaging over the set of such differences estimates the treatment effect. So although individuals have not actually been randomly allocated to treatments, the fact that overall distributions of covariates differ between the groups is “ignorable”, given matching on the propensity score.

A Propensity Score analysis commences with estimation, by probit or logit, of a treatment assignment equation, where all known covariates affecting assignment and response are included as explanatory variables and the observed “dependent” variable is $D=1$ for an individual in the treatment group and $D=0$ for someone in the control group. Then propensity scores are calculated for all individuals and some matching process is implemented. The most commonly employed is stratification of the

propensity score distribution by quintiles or sextiles - the "binning" procedure. Then the distributions of covariates for treatment and control within each subclass are compared and, if they still differ appreciably, the assignment equation is further developed. For example, if a particular covariate still differs between groups within subclasses, the assignment model could be modified by trying powers of the covariate and its interactions with other variables. It is important to search for the best model for participation. When a satisfactory model is arrived at, the treatment v control effect on the response variable within subclass i is just the difference in means $\bar{y}_{1i} - \bar{y}_{0i}$, if the response is continuous, or a difference in proportions $\hat{p}_{1i} - \hat{p}_{0i}$, if the response variable is qualitative. Then the overall measure of treatment effect can be taken as simply

$$\frac{1}{s} \sum (\bar{y}_{1i} - \bar{y}_{0i}), \quad (9)$$

where s is the number of bins, or strata, in which there are both treatment and control units and the summation is over these strata. The standard error is

$$\frac{1}{s} \sqrt{\sum \left\{ \frac{\hat{\sigma}_{1i}^2}{n_{1i}} + \frac{\hat{\sigma}_{2i}^2}{n_{2i}} \right\}}, \quad (10)$$

where the $\hat{\sigma}_{1i}^2$ and $\hat{\sigma}_{0i}^2$ are the within group and within stratum i variances among the n_{1i} treated individuals and the n_{0i} controls. For example,

$$\hat{\sigma}_{1i}^2 = \frac{1}{n_{1i} - 1} \sum_j (y_{1ij} - \bar{y}_{1i})^2.$$

If the response variable is qualitative, formulae (9) and (10) become

$$\frac{1}{s} \sum (\hat{p}_{1i} - \hat{p}_{0i}) \quad \text{and} \quad \frac{1}{s} \sqrt{\sum \left\{ \frac{\hat{p}_{1i} \hat{q}_{1i}}{n_{1i}} + \frac{\hat{p}_{0i} \hat{q}_{0i}}{n_{2i}} \right\}}, \quad \text{with } q = 1 - p.$$

The choice of (9) as an estimate needs some discussion. If the treatment versus control comparison could be presumed the *same* within all strata (that is, training has the same consequences for a low propensity individual as for a high one) then (9) is not the *best* estimate, although it is an unbiased one. The best estimate, in a minimum variance sense, would weight the within strata estimates inversely as their variances (so giving most weight to the most precise estimate). But if we want to allow for varying treatment effects across strata, we need to weight each stratum contrast in proportion to the

stratum's fraction of the population, which means equally, given quintiles or sextiles. The constant treatment effect assumption is implicit in the econometric approach of Section 2, but is not essential for the Propensity Score approach. However, it may well be plausible in some cases and can permit more precise estimates and more powerful tests of treatment effects.

It is possible that the Propensity Scores approach could fail to achieve a comparison of treatment with control. The difference within stratum i , $\bar{y}_{1i} - \bar{y}_{0i}$, obviously presupposes that there are some treated and control individuals present. If there are no representatives of one group, that stratum does not contribute to the comparison. If no stratum contains representatives of both groups, the approach fails, because there is no overlap in the propensity scores. The interpretation is that the characteristics (as measured by covariates) of the two groups are so dissimilar that no meaningful comparison is possible. This is not a disadvantage of the Propensity Scores method relative to the econometric modelling approach. The data deficiencies would feed into and undermine the regression analyses, although the cause of the problem might not seem at all obvious. One of the virtues of the Propensity Scores approach is that it reveals just how much of the data truly provide information on the comparison. Dehejia and Wahba (1998), who reanalysed Lalonde's (1986) data, emphasise this point. Rubin (1997), writing in the context of drawing deductions from health care databases, has stressed that the first use of propensity scores should be to decide whether a question of interest can be legitimately addressed to the database at all.

Some other advantages of the method deserve mention at this point. The approach is nonparametric as regards the response variable and is very sparing on assumptions. Nothing has been specified about the actual relationships of the response to the covariates, so avoiding the accumulation of biases due to the combination of model misspecifications and unbalanced covariates, which can have serious consequences for the regression modelling approach (for example, see Rubin, 1997). The estimation and testing of the treatment v control difference will usually be more efficient because the standard error will be smaller. Several factors contribute to this, but the most important is that multicollinearity induced variance inflation, so often associated with multiple regression on a lot of correlated variables, is being avoided. The reduction from multidimensional covariates to a unidimensional propensity score also makes results easier to summarise and understand. This point might seem trivial, but it recurs frequently in the literature comparing (for example, Rubin, 1997; Oberchain and Melfi, 2000; Perkins et al, 2000) the Propensity Scores approach with multiple regression/econometric/Heckman methods.

At this point it should be said that there are other matching procedures besides stratifying into quintiles or sextiles. More subclasses could be employed, for example by stratifying into deciles. Or bins could be based on equal ranges, rather than frequencies, of propensity scores as in Tu, et al (2000). Leaving stratification entirely, each treatment individual could be matched with the control individual with the closest propensity score value, or matched to a group of “close” control cases. Decisions on what constitutes “close” can be based on “callipers” – pre-selected ranges. Although these procedures sometimes have advantages, most applications of the Propensity Score methodology use “binning” because Cochran (1968) showed that stratification into quintiles usually removes 90% of the bias due to differing covariate distributions between the treatment and control groups. It should also be said that multiple regression, or econometric, modelling can have a refining role, if applied within strata following matching on the propensity score (see Drake and Fisher, 1995, or Rubin, 1997). This is because model misspecifications cause minimal biases if the covariate distributions are similar for the treatment and control groups. While it is probably obvious, the y variables in (9) could be replaced by differences of post- and pre-treatment (and control) values if the earlier measurements exist. This would give a matched or Propensity Score “difference of differences analysis”, very analogous to common practice in randomised experiments.¹²

When there are more than two groups comparisons are made pairwise, with separate derivation of propensity scores for each pair. The need for this can be appreciated by extending (following Rubin, 1997) the smoking example to three groups – non-smokers, cigarette smokers and pipe smokers – and two covariates – age and social class. As before a mortality measure is the response variable. Suppose non-smokers and cigarette smokers have the same age distributions, but unequal, though overlapping, social class membership. Suppose non-smokers and pipe smokers have the same distributions by social class, but unequal, though overlapping, age distributions. Then for the non-smokers v cigarette smokers comparison the propensity score matching should “balance” social class differences, while for non-smokers v pipe smokers it should “balance” age differences. Clearly separate derivations of propensity scores are appropriate.

¹² For example, in randomised animal growth experiments, the difference between mean final weights of treatment and control groups is a valid estimate of treatment effect, but the difference in mean weight gains (assuming initial weights were recorded prior to treatment application) is the preferred measure, because it reduces random error.

Matching on the propensity score balances over observed covariates. What if an important covariate – one that has a substantial effect on the response *and* whose distribution differs considerably between treatment and control groups – has not even been observed? If it, *w* say, is uncorrelated with the observed covariates, matching on the propensity scores can no longer be assumed to have adequately simulated a randomisation situation. However, if it is correlated with them, matching on the propensity scores based on the observed covariates will also, at least partially, balance the unobserved covariate effect. This is because such matching produces treatment and comparison groups with the same distribution of *x* variables. They would balance the unobserved variable most effectively if *w* had a well fitting regression on the observed variables with a disturbance uncorrelated with *y*. If the regression explained enough of the variation in *w*, they would still balance effectively even if the disturbance was not uncorrelated with *y*. For this reason the Propensity Score approach stresses searching for and examining the maximum number of attainable covariates. The Heckman correction procedure given by (5) and (6) also assumes *w* has a linear regression on the observed covariates, though with the additional assumption of bivariate normality (and, some critics would say, without the stress on searching). The normality assumption would be invalidated by the existence of a non-random unrecorded covariate. In reality, both the Propensity Score and the econometric modelling/Heckman correction approaches assume the observed covariates are sufficient to work with, but the difference between the approaches can be very clearly seen by imagining that *w*, “ability” say, becomes observable. Now the econometric procedure just involves estimation of equation (3) with *w* added as an extra variable. There is no estimation of a probit or logit participation equation, because the selection bias was presumed to operate only through the formerly unobserved variable. But the Propensity Score method will, of course, retain the treatment assignment equation, just adding *w* to the variables in it, on the supposition that all variables can contribute to selection bias.

However, if an unobserved, covariate is not well explained by observed covariates and has a substantial effect on *y*, the “strongly ignorable assumption” – that treatment assignment is effectively random within strata of the propensity score – does not hold and progress would depend on other assumptions. It is possible to employ IV methods within a Propensity Score framework (for example, Angrist, 1995) provided truly valid instruments can be identified, but this an issue about which there are considerable disagreements in the recent technical literature. Before returning to this topic it is

desirable to review the measure of agreement that has been reached about the Propensity Score approach in this literature as well as the caveats that have been expressed and the unresolved issues.

IV RECENT DEBATES IN THE LITERATURE

A quick reading of recent literature can give a first impression of substantial, even acrimonious, disagreement between authors with the opposing sides (largely) consisting of econometricians on the one hand and “mainline” statisticians on the other. But at least some of this is due to rival claims about who deserves the credit for the formulation of key concepts and to criticisms of each others’ originality and presentation of approaches. In fact, there is more or less general agreement about the desirability of “balancing” covariate distributions through matching and acceptance that, in most if not all cases, propensity scores, or closely related statistics, should be employed. It is now generally appreciated that, whatever about selection bias arising from an unobserved characteristic, considerably larger biases can follow from failure to correct adequately for very different covariate distributions between treatment and control.

Heckman, Ichimura, Smith and Todd (1996) said that their data analysis “appeared to provide a strong endorsement for matching on the propensity score”. They stressed that incomparable individuals should not be used in contrasting the treatment and control (in a “binning” context, strata containing individuals from only one group should be discarded), or selection bias would result. But their comment that this “essential” requirement was “hitherto unnoticed” is unfair to many previous authors outside econometrics. The data analysis in Heckman, Ichimura and Todd (1997) decomposed bias into three components: that due to non-overlapping support (the comparing of incomparables), that due to different distributions of covariates within groups and finally that due to selection on unobservables. They said that this last component “called selection bias in econometrics” had been emphasised in the previous econometric literature, but was actually smaller in magnitude than the others and they concluded “Simple balancing of observables . . . goes a long way towards . . . effective evaluation . . .”

This is not to say these authors were uncritical of Propensity Score methodology. They believed that bias due to selection on unobservables, even if it was smaller than the other components, required attention, and this topic will be returned to shortly. They also made other criticisms and comments that are perhaps not of great practical import in most situations. Heckman et al (1996) said that because

much of the data might be useless (the individuals having no matches in the other group) a Propensity Score analysis of survey data could be greatly inferior to a true randomised experiment of equal size. The latter would also handle selection on unobservables, since randomisation “balances” over both observed and unobserved covariates. The points¹³ are certainly true if the researcher has that choice, but are of no relevance if the choice is between Propensity Score and the regression/Heckman correction approach of Section II. Unlike our implicit assumption that the conceptual comparison of interest is the mean effect of treatment on the population from which both treatment and control subjects were drawn versus the mean “effect” of the control on that population, the papers mentioned take the conceptual comparison as the effect of treatment on the treated. That is, the control group simulate how the treatment group would have performed without the treatment. This makes no difference to the mechanics of matching etc, but may matter for interpretation and, the authors argue, may be less demanding as regards assumptions.

Both papers raise a point followed up in Heckman, Ichimura and Todd (1998), and Heckman, Ichimura, Smith and Todd (1998) – the desirability of matching need not imply that propensity scores are the only, or indeed the best, way to achieve the matching. Rosenbaum and Rubin (1983) proved that matching on propensity scores would balance all covariates, but their proof assumed the propensity scores known exactly, rather than estimated. The cited papers show that matching on estimated propensity scores may not be as efficient (in the sense of attaining derived asymptotic bounds to precision) as matching on all the covariates. But the force of the argument fades if such matching is possible only in the rare cases where data sets are huge, relative to the number of covariates. The original motivation for the propensity score rather than Cochran’s (1965, 1968) matching on covariates was practicality rather than efficiency. However, there is more to be said than this. In general, estimated propensity scores could differ from true values either because the estimating equation was misspecified or because sampling errors could be appreciable with a finite data sample. With such very large data sets, the second possibility can be neglected and the first surmounted by employing a nonparametric method. If covariates are categorical, there will be grouped observations from which participation probabilities can be directly obtained, while Kernel estimation can achieve the same result for continuous covariates.

Somewhat similarly, although in a purely theoretical paper, Hahn (1998) argues that an efficient

¹³ These views are perhaps difficult to reconcile with the Heckman and Smith (1996) view of randomisation as a “misleading paradigm” for the social sciences.

estimator of average treatment effect is attainable without matching on the propensity score (even if it is known) by “nonparametric imputation” and “various nonparametric regression techniques”. As regards known propensity score, while it is not at all implausible that there may be another way to obtain an efficient estimator, it is hard to see why it should not be done the easy way. In general though, the previous comments about limitation to cases of copious data must still apply. Informative nonparametric methods usually demand a lot of data, although the details of data requirements are not discussed in the paper. Finally, if Hahn’s claim that matching on the propensity score is sometimes worse than not than matching at all is to be taken seriously, it requires more evidence than saying the data may have been obtained from a properly randomised experiment. If researchers know they have a randomised experiment, they will of course analyse it as such.

We think the implications of having to estimate propensity scores are more realistically assessed in the context of data sets of a size requiring parametric estimation via probit or logit models and with emphasis on finite sample rather than asymptotic criteria. Somewhat contrary to the impression given by the recent econometric literature, the consequences of such estimation have been considered in the biometrical/statistical literature. For example, Drake (1993) employed simulation studies to look at effects of various forms of misspecifications of the propensity score estimating equation, including omitted variables, incorrect functional form, etc. She found that many misspecifications did not matter very much, provided all independently important covariates had been included in the estimating equation. It is worth noting too, that estimated propensity scores do not have to prove inferior in practice to known true values. Rosenbaum (1987, section 2.3) explains why estimated propensity scores can often *outperform* the true values.

Other points raised in the econometric literature are either not really contentious or are rather peripheral. For example, Heckman et al (1998) suggest that, if a balancing score must be used (rather than balancing on all the covariates), the inverse Mill’s ratio (Heckman’s lambda) could replace the propensity score. Since one is a monotonic function of the other, this should make little, if any, difference. The statistical literature has itself considered balancing scores other than defined as the probability of inclusion in the treatment group. The discriminant function (for example, Rosenbaum and Rubin, 1985) and classification tree methods have been employed. The former can also be shown to be a monotonic function of the propensity score if the discriminant function is that appropriate to discrimination between two multivariate normal distributions with the same covariance matrix,

although that is unlikely given the frequency of categorical variables in typical data sets. The relevance of other measures of difference in response between groups, besides mean difference, comes up in some papers (see also: Heckman, Smith and Clements, 1997; Imbens and Rubin, 1997). The treatment effect may be thought of as “random” (varying over individuals) rather than “fixed” leading to different distributions (not just different means) of the response variable in both groups. This is a topic on which much more may well appear in the future. Random treatment effects (sometimes called components of variance models) have a long history in the analysis of experiments and have also been frequently employed in econometrics, especially in the analysis of combined cross-sectional and time series data. However, the vast majority of applied papers where adjustment for selection bias has featured, whether by propensity score or econometric modelling (including the Irish papers cited in Section I) have taken the fixed treatment effects context as appropriate and the selection bias software routines in econometric packages invariably assume it.

Returning to the topic of selection bias due to unobserved covariates, Heckman et al (1997) and Heckman et al (1998) used large data sets¹⁴ to test for the existence of residual selection bias (due to selection on unobservables) following matching for observables. They found evidence of such effects in their data and although the magnitudes were small relative to the biases eliminated through matching, they were still appreciable relative to the size of treatment effect. (Care is required with this kind of argument – if the treatment effect is zero, or near zero, any bias is appreciable relative to it.) They considered two approaches to the correction of this residual selection bias. One utilised the fact that they possessed repeated measurements over time on the response variable (longitudinal data) and they constructed a variant of the “difference of differences” estimate of treatment effect, which was successful in eliminating this residual bias. Now as pointed out in the previous section, use of such an estimate within matched strata is in the tradition of experimental analysis and fits naturally into the Propensity Score methodology. Even if there was no bias due to selection on unobservables, the use of the procedure would be justifiable in terms of reducing random variation and hence increasing precision. The reason it frequently does not feature in biometrical or epidemiological applications of the Propensity Score method is because the post treatment response measure is often all that is available.

¹⁴ They actually had a true randomised control group as well as several “observational data” control groups that required the matching approach.

Without pre treatment response measures, correction for selection on unobservables has to depend on behavioural assumptions, which could perhaps be wrong and are often not really testable. The possibility of using instrumental variables was mentioned at the end of the last section, but on this matter there is a definite difference of opinion between econometricians and “mainline” statisticians. Heckman, Ichimura and Todd (1998) argue that knowledge of “zero exclusions” should be integrated into analyses and show how to do so. This may correct for selection bias arising from unobservables, when such bias exists, but even if it does not, integration can (they argue) be justified because using all the available information must lead to a more precise estimator of treatment effect. Statisticians seem to have less faith in the reliability of the assumptions made by econometricians. They do not deny the possibility of selection bias due to failure to balance unobservable covariates and Rosebaum (1984, 1989) has suggested relevant tests, while Rubin (1997) has suggested sensitivity analyses on the lines of Rosenbaum and Rubin (1983b). However, they believe that substantial biases can be introduced by inexact behavioural (including zero restrictions) assumptions. Little and Rubin (1999) describe models incorporating them as “highly sensitive to minor deviations from assumptions” and Oberchain and Melfi (2000) say their experience is that “models using . . . instrumental variable analyses quickly become . . . frustratingly sensitive to the validity of . . . underlying assumptions”. Little and Rubin argue that such introduced biases can easily exceed the residual biases from unbalanced unobservables, especially given a comprehensive search for observable covariates. Supporting evidence is cited of the effects of departures from the “ignorability” condition on the (conceptually highly related) issue of dealing with missing values in surveys (for example, Rubin, Stern and Vehovar, 1995).

So on this issue of how to deal with selection on unobservables there is still disagreement in the literature, but it should not obscure how much agreement there is on other matters. For completeness, a quite different approach to correcting for selection bias – placing bounds upon it – will be briefly described. As with matching, this approach goes back to Cochran (1953), who employed it in coping with non-response in surveys. An example can illustrate the original form of the approach. If, say, in estimating the proportion of voters for party A, one fifth of the sample refuse to respond, but among respondents half indicate preference for A, the bounding procedure would operate as follows. The “true” proportion is four fifths the proportion among respondents plus one fifth the proportion among non-respondents. The latter must be between zero and one. So an estimate of the true proportion lies between two fifths minus zero = two fifths and two fifths plus one fifth = three fifths. The approach

has been developed and applied to selection bias by Manski (1990) and Horowitz and Manski (1998). However, the bounds are often found to be too widely separated to be really useful, and we will not pursue the approach further.

V EXPOSITORY ANALYSIS OF IRISH TRAINING PROGRAMME DATA

O'Connell and McGinnity (1997) examined the effectiveness of Irish educational training and employment schemes conducted by the State training Authority (FAS) and the Department of Education. In mid-1994 they interviewed a large sample (4600) of individuals who had exited training courses between April and July 1992 and ascertained their pre and post training labour market experiences, as well as information on their education, family backgrounds and social circumstances. The courses fell into several categories as regards the type of training and only one category – general training – will be considered here. This is adequate given the expository context of this paper, although the same approaches could be applied to the other course categories. General training courses provided instruction in a range of basic skills and were mainly intended for people with relatively poor educational qualifications.

O'Connell and McGinnity constructed a control group by selecting suitable people from the Economic and Social Research Institute's long running School Leavers Survey. This survey takes annual samples of school leavers and follows the cohorts over subsequent years. The criteria for selection were that individuals had left school between 1990 and 1992, were unemployed and in the labour market at the same time as trainees were exiting programmes, and had not participated in training courses themselves. They were also interviewed in mid-1994. Clearly these people were relatively young and, to avoid an obvious source of comparison bias, O'Connell and McGinnity excluded all trainees aged over 23 from the analysis. Nonetheless there were considerable differences between the treatment and control groups in some other possibly relevant characteristics. Table 1 compares the training and control groups in terms of these characteristics, or covariates, showing mean values for continuous covariates and proportions for categorical characteristics and indicating statistically significant differences between groups. Note, in particular, the considerable differences in educational attainments. On average at least, the control group are more advantaged in terms of education and other socio-economic characteristics.

Table 1: Analysis of Covariates by Treatment and Control

	Treatment	Control	Tests for Difference	
	Mean	Mean	<i>T-Test</i>	
Age	18.46	18.44	0.21	
Household Size	5.08	4.85	1.59	
Duration of Unemployment	4.54	3.27	3.20	**
	%	%	<i>Chi-Square</i>	
Female	43.77	47.97	1.35	
No Qualifications	33.91	4.94	79.20	***
Junior Certificate	37.36	27.16	8.55	**
Leaving Certificate	25.77	65.84	131.96	***
Third Level Education	2.96	2.06	0.57	
Respondent had never worked	67.92	90.24	47.94	***
Father in Employment	42.82	55.00	11.04 **	
Mother in Employment	16.11	14.58	0.32	
Father Employed at School Stage	54.18	65.13	8.93	**
Mother Employed at School Stage	7.76	6.22	0.64	
Fathers Social Class:				
Professional	11.64	16.45	6.03	*
Non-manual Skilled	44.25	47.19		
Semi-unskilled/Manual	44.11	36.36		

Significance of P at the following levels: * $p < .05$, ** $p < .01$, *** $p < .001$

The response variable of interest was taken to be employment status eighteen months after completion of training. In line with the exposition of Section II, an equation of the form (3) (but a probit model because of the binary response variable) was estimated containing a dummy variable to contrast training and control groups and a range of covariates to try to cancel out the other differences between groups. Several equation specifications, differing in the actual set of covariates included, were tried out before that shown in Table 2 was chosen.¹⁵ Withholding some (non-significant) covariates from the model is essential, as was explained in Section II, if Heckman style sample selection bias testing is to be conducted with a binary response variable.

Table 2: Probit Model of Employment after 18 months: General Training v Control Group

	Coefficient	Standard Error	t-ratio	P value
Constant	0.077	0.510	0.150	0.880
General Training	0.024	0.105	0.229	0.819
Female	-0.176	0.085	-2.072	0.038
Age	-0.043	0.031	-1.403	0.161
Junior Cert.	0.565	0.118	4.788	0.000
Leaving Cert.	1.028	0.142	7.253	0.000
Unemployment Duration (pre-programme)	-0.023	0.008	-3.079	0.002
Log Likelihood	-631.4	Chi ² 98.3	No cases	1011

There is no statistically significant effect of general training, although some covariates do clearly impact on the likelihood of job attainment. Continuing to a test for sample selection bias, our model

¹⁵ This analysis is not identical to that presented in O'Connell and McGinnity (1997) because their analysis was conducted simultaneously for all the categories of training and employment schemes.

for the training group is

$$y_j^* = a + \sum_{k=1}^{k=p} c_k x_{kj} + u_j, \quad (11)$$

where y^* denotes the latent variable underlying the observed binary response variable y and the training effect is considered contained in the intercept, a . The participation equation, which applies to the control group as well as to the training group is

$$D_j^* = g + \sum_{k=1}^{k=s} h_k x_{kj} + v_j, \quad (12)$$

where D^* is the latent variable underlying the binary variable D , which defines membership of the training or control groups. Sample selection bias exists if a component of v , the disturbance term in (12), is correlated with u , the disturbance term in (11), for the common training observations. It can be tested for by estimating (by maximum likelihood) this correlation and seeing if it is significantly different from zero. The procedure can be repeated with the control group replacing the treatment group in (11) to test for selection bias (possibly of different magnitude) in the control group.¹⁶ The LIMDEP package (Green, 1991) can perform the required computations. The results are in Table 3.

Table 3: Testing for Sample Selection Bias for the Employment effects of General Training

	ρ	Std. Error	t-ratio	Signif.
General Training	- 0.289	0.264	- 1.095	0.273
Control group	- 0.281	0.371	- 0.756	0.450

The estimates of correlations are not significant and so there seems no reason to modify the conclusion, as drawn from Table 2, that general training does not improve an individual's prospect of gaining employment.

Turning now to the Propensity Score approach, we sought logistic models¹⁷ for the probability of participation using combinations of the variables listed in Table 1. Subsequent modifications introduced other variables, generated as powers or interactions of existing variables. The choice of model was partly on standard statistical criteria – the significance of the coefficients (jointly as well as

¹⁶ The test could have been applied assuming the same covariate coefficients in the treatment and control groups, in line with equation (3), but the somewhat more general approach adopted here has been recommended in the literature (for example, Maddala, 1983, p. 261).

¹⁷ SPSS was employed because some of the Propensity Score procedures are easily represented graphically with that package. However, it was then preferable to use a Logistic rather than a Probit formulation, because SPSS provides only limited options for the latter. The propensity scores themselves should be very similar whether produced from Logistic or Probit.

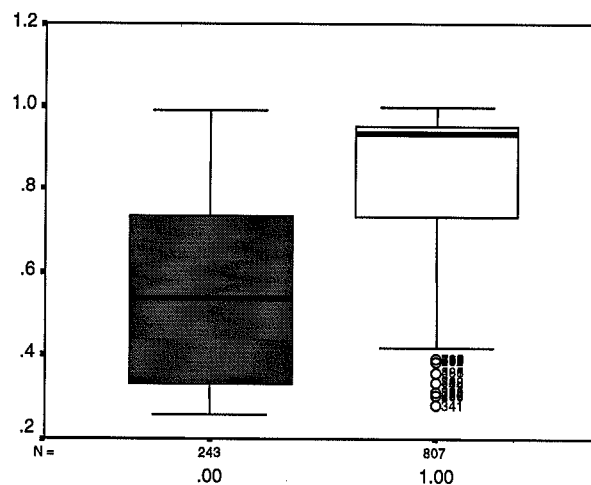
individually), the fit of the model and the retention of observations (some covariates were much more prone to missing values than others). More importantly, however, the procedure outlined in Section III was also followed. Models were assessed on how well they achieved within stratum balance of the propensity score and of the covariates. Because there was a considerable degree of collinearity between covariates, many models were quite similar in their predictive ability. Although parsimony of model covariates is *not* particularly desirable in modelling the propensity score, for expository purposes we focus on the reasonably simple model shown in Table 4.

Table 4: Logistic Model for Participation in Training

Covariate	Coefficient	SE	Wald Test	P value
Constant	37.98	11.96	10.1	.0015
Female v Male	.14	.18	.6	.4306
No Qualifications	3.65	.37	96.3	.0000
Junior Cert.	1.80	.22	66.54	.0000
Never Worked	-1.67	.25	44.65	.0000
Age	-4.36	1.28	11.58	.0007
Age squared	.128	.034	13.82	.0002
Father Skilled	-.004	.211	.00	.98
Father Semi or Unskilled.	-.118	.215	.30	.58
Log Likelihood	-422.5		Chi ²	291
R ² (Nagelkerke)	.366		No. of cases	1050

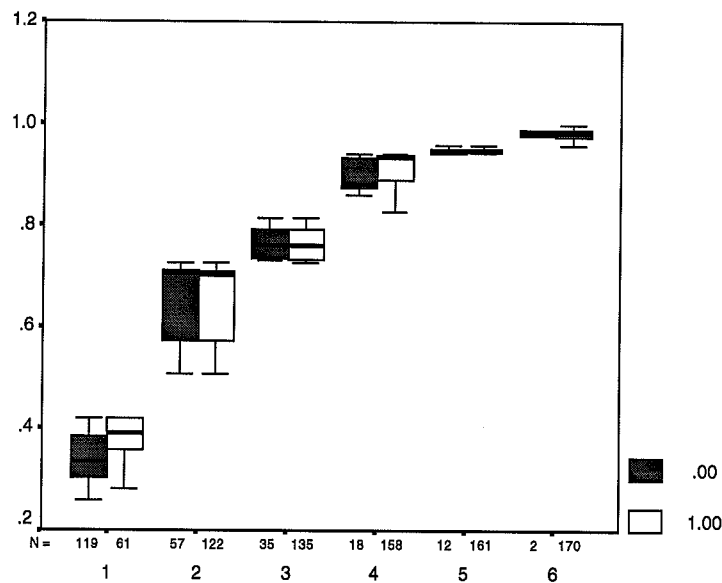
The high significance of some covariates (the two dummy variables for education, for example) was to be expected given the differences shown in Table 1, but age is more surprising and must follow from different distributions of age (in spite of near equal means) in the two groups. Box plots of the propensity scores produced by this model, categorised by treatment and control, are shown in Figure 1.

Figure 1: Propensity scores for Control v General Training



Clearly there is no balance between groups in the distributions of propensity scores, with the training group's values much higher overall. The result of stratifying the distribution into sextiles and comparing distributions within the strata (bins) is shown in Figure 2.

Figure 2: Within Strata Comparisons of Propensity Scores



Clearly the discrepancies between groups are much reduced within bins. Only in the lowest sextile is there still a substantial difference (and statistically highly significant) between treatment and control (and perhaps the Bin 4 difference is greater than desirable, though not statistically significant).

However, the balance has been achieved at the price of different distributions across bins. Table 5 shows that most of the treatment group are in the upper sextiles and most of the control group in the lower sextiles. There are only 2 individuals in the control group in bin 6.

Table 5: Frequencies across Sextiles and Mean Propensity Scores (PS)

	Bin 1	Bin 2	Bin 3	Bin 4	Bin 5	Bin 6	Overall
Frequencies: Treatment	61	119	133	156	160	167	796
Control	119	57	35	18	12	2	243
PS mean: Treatment	.383	.665	.764	.916	.948	.981	.832
Control	.346	.661	.764	.900	.949	.983	.556
PS mean for bin	.359	.664	.764	.914	.948	.981	.769

As discussed in previous sections, balancing over propensity scores should balance over the included covariates and this can be checked out. However, this obviously has to be interpreted in the light of the frequency distribution between strata. It would be quite probable that the control group in bin 6 would show either 0% or 100% for a binary characteristic, perhaps suggesting even less balance after

stratification than before. But, of course, this would not be statistically significant.

It is probably more interesting to look at the degree of balance achieved over a covariate that has *not* been included in the model. Table 6 shows the between and within strata variation of the variable "duration of unemployment" (pre training).

Table 6: Duration of (pre training) Unemployment in months

	Bin 1	Bin 2	Bin 3	Bin 4	Bin 5	Bin 6	Overall
Duration UE : Treatment	2.45	3.48	3.66	4.24	5.30	6.33	4.54
Control	1.78	3.84	3.86	7.39	7.50	nc	3.27
Statistical significance	ns	ns	ns	≈ *	ns	ns	**
Bin Means for duration	1.99	3.60	3.70	4.57	5.46	6.30	4.24

Significance of P at the following levels: *p<.05, **p<.01, ***p<.001

ns = not significant = not calculated

This covariate did differ very significantly between treatment and control and in fact it would have improved the propensity score model somewhat if added in. But clearly most of its influence is being captured by included variables that are correlated with it and indeed the regular increase in duration of unemployment with propensity score, as shown by the bin means across strata, makes this intuitively obvious. Only for bin 4 did the within stratum difference approach 5% significance. (Of course, frequencies of the control in bins 4 and 5 are low). This illustrates the discussion in previous sections about approaches to, and consequences of, an unobserved covariate. Socio-economic variables are often quite correlated with each other and so if substantial efforts have been made to take account of relevant covariates in estimating the propensity score, a considerable amount of balance may also be achieved over an unobserved variable.

Effect of Training on Subsequent Employment

With only 2 individuals in the control group in bin 6, it is clear this stratum itself cannot provide useful information about the treatment effect on the dependent variable, since the latter is the binary variable - had/had not a job 18 months after training. However, it seems a pity to have to discard the 167 treatment values, especially since Table 5 and Figure 2 show that overall bin means are converging in the upper sextiles. The difference in mean propensity score between bins 1 and 2 is .3, while between bins 5 and 6 it is only .03. So in Table 7, which compares the effects of treatment and control on the proportions in employment eighteen months after the completion of the training courses, bins 5 and 6 have been combined. Even so, there are only 14 values for the control, so the within combined stratum comparison will not be at all precise. In the table "Overall" means the comparison *without* any matching on propensity score.

Table 7: Proportions Employed 18 months post-training

	Bin 1	Bin 2	Bin 3	Bin 4	Bins 5+6	Overall
Proportions : Treatment	.49	.50	.43	.31	.30	.37
Control	.63	.39	.43	.17	.35	.49
Statistical significance	≈ *	ns	ns	ns	ns	***
Bin Means for Proportions	.58	.47	.43	.30	.30	.40

Significance of P at the following levels: * $p < .05$, ** $p < .01$, *** $p < .001$
 ns = not significant = not calculated

Note first from the bin means that the proportions getting employment decline with increasing propensity score. That is perfectly compatible with the facts, remarked earlier, that the control group were more advantaged in terms of educational qualifications etc., so that these characteristics would be negatively related to the probability of being in the training group. But these characteristics would help individuals get employment. So the overall (that is, ignoring the propensity scores) finding that the control group had a highly significant advantage in obtaining employment (49% to 37%), does not demonstrate the failure of training. The important thing is how treatment and control groups compare at similar propensity scores. Within the bins none of the differences are statistically significant at the 5% level, but that for bin 1, which favours the control, approaches it and is significant at the 10% level. However, we recollect from Figure 2 and Table 5 that balance was not attained within this bin in that the propensity score was substantially lower for the control. So we refine the binning procedure for the lowest sextile alone by splitting it into halves. Figure 3 shows the overall comparison of propensity scores for the sextile and Figure 4 shows the result of halving.

Figure 3: Con. v Gen. Train. - Bin 1

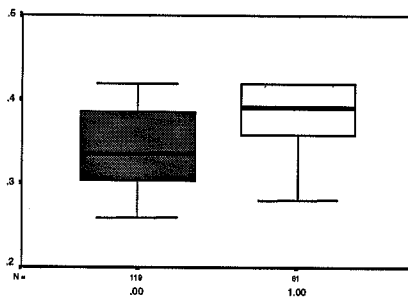
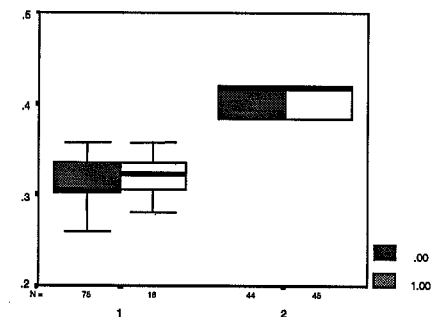


Figure 4: Con. v Gen. Train. - Bin 1 Split



Balance on the propensity score is now effectively achieved. In one substratum the proportion employed is now very slightly higher for the treatment group, while in the other it is lower for the treatment group, although not sufficiently so to be statistically significant.

The estimation formulae in Section 3 assumed equal proportions of the data in each stratum and

need modification now. For the average treatment effect the formula is now

$$\frac{1}{s_i} \sum (\hat{p}_{1i} - \hat{p}_{0i}),$$

where the within stratum difference in proportions employed is weighted by the stratum's proportion of all individuals. So for the halves of the lowest sextile the s 's are 12, for the 2nd, 3rd, and 4th sextiles they equal 6, and for the combined 5th and 6th, s equals 3. The result is an overall small (.009) difference in favour of the training group, but since its standard error turns out to be .041, it is far from statistically significant.

Besides lack of balance there may be another reason to consider bin 1 separately from the others. Table 5 showed that the biggest jump in propensity scores occurs between bins 1 and 2, with the mean for the latter .66. Perhaps we should accept that general training may not be particularly beneficial to the relatively well educated and advantaged young people who constitute the population of bin 1. Could training be beneficial for the less advantaged (those with higher propensity scores)? Although treatment was not significantly better than control within any one bin, the criteria were approaching critical values in bins 2 and 4. A method of assessing if a set of r non-significant tests attain joint significance (Fisher, 1932) is to insert the P_i values from the tests into the formula

$$\chi_{2r}^2 = -2 \sum \log(P_i).$$

For bins 2 to 5+6 this gives a value of 11.7 which, for a chi-squared with 8 degrees of freedom, corresponds to a P value of .15 and so is short even of 10% significance. So we cannot assert a definite advantage for training.

The Propensity Score approach has not led to different conclusions than the "classical" analysis of Tables 2 and 3 although it has made far fewer assumptions and, we think, has probed deeper into the data structure. The general training v control comparison was just chosen for expository purposes and it will be interesting to see if the approach confirms or conflicts with other reported findings from the application of econometric modelling and Heckman correction to Irish data on labour market interventions. We hope, however, that we have adequately conveyed the key ideas of the Propensity Score approach, outlined its scope and illustrated its application.

REFERENCES

- ANGRIST, J. D. AND A. KRUEGER, 1991. "Does Compulsory School Attendance Affect Schooling and Earnings", *Quarterly Journal of Economics*, Vol. 106, pp. 979-1014.
- ANGRIST, J. D., 1995. "Conditioning on the Probability of Selection to Control Selection Bias", *Technical Working Paper 181*, Washington: NBER.
- ANGRIST, J. D., G. W. IMBENS AND D.B. RUBIN, 1996a. "Identification of Causal Effects Using Instrumental Variables" (with discussion), *Journal of the American Statistical Association*, Vol. 91, pp. 444-472.
- ANGRIST, J. D., G. W. IMBENS AND D.B. RUBIN, 1996b. "Rejoinder to Heckman", *Journal of the American Statistical Association*, Vol. 91, pp. 468-472.
- BREEN, R., 1986. *Subject Availability and Student Performance in the Senior Cycle of Irish Post-Primary Schools*, ESRI General Research Series Paper No. 129, Dublin: ESRI.
- BREEN, R. and B. HALPIN, 1988. *Self-employment and the Unemployed*, ESRI General Research Series Paper No. 140, Dublin: ESRI.
- BREEN, R. and B. HALPIN, 1989. *Subsidising Jobs: An Evaluation of the Employment Incentive Scheme*, ESRI General Research Series Paper No. 144, Dublin: ESRI.
- BREEN, R., 1991. "Assessing the Effectiveness of Training and Temporary Employment Schemes: Some Results from the Youth Labour Market", *The Economic and Social Review*, Vol. 22, pp. 177-198.
- BREEN, R., D. HANNAN and R. O'LEARY, 1995. "Returns to Education: Taking Account of Employers Perceptions' and Use of Educational Credentials", *European Sociological Review*, Vol. 11, pp. 59-73.
- COCHRAN, W. G., 1953. *Sampling Techniques*, New York: Wiley.
- COCHRAN, W. G., 1965. "The Planning of Observational Studies of Human Populations" (with discussion), *Journal of the Royal Statistical Society A*, Vol. 128, pp. 234-255.
- COCHRAN, W. G., 1968. "The Effectiveness of Adjustment by Subclassification in Removing Bias in Observational Studies", *Biometrics*, Vol. 24, pp. 205-213.
- D'AGOSTINO, R. B., 1998. "Propensity Score Methods for Bias Reduction in the Comparison of a Treatment to a Non-randomised Control Group", *Statistics in Medicine*, Vol. 17, 2265-2281.
- DEHEJIA, R. H. and S. WAHBA, 1999. "Causal Effects in Non-experimental Studies: Reevaluating the Evaluation of Training Programmes", *Journal of the American Statistical Association*, Vol. 94, pp. 1053-1062.
- DRAKE, C., 1993. "Effects of Misspecification of the Propensity Score on Estimators of Treatment Effect", *Biometrics*, Vol. 49, 1231-1236.
- DRAKE, C. and L. FISHER, 1995. "Prognostic Models and the Propensity Score", *International Journal of Epidemiology*, Vol. 24, pp.183-187.
- EISSA, N. and J. B. LIEBMAN, 1996. "Labour Supply Response to the Earned Income Tax Credit", *Quarterly Journal of Economics*, Vol. 111, pp. 605-637.
- FISHER, R. A., 1932. *Statistical Methods for Research Workers*, 4th edn, London: Oliver and Boyd.
- GOLDBERGER, A. S., 1983. "Abnormal Selection Bias", in Karlin, S., T. Amemiya and L. A. Goodman (eds.), *Studies in Econometrics, Time Series and Multivariate Statistics*, New York: Academic Press, pp. 67-84.
- GREEN, W. H., 1991. *LIMDEP Version 6.0*, New York: Econometric Software.
- HAHN, J., 1998. "On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects", *Econometrica*, Vol. 66, pp. 315-331.
- HECKMAN, J. J., 1976. "The Common Structure of Statistical Models of Truncation, Sample Selection Bias and Limited Dependent Variables and a Simple Estimator of such Models", *Annals of Economic and Social Measurement*, Vol. 5, pp. 475-492.
- HECKMAN, J. J., 1979. "Sample Selection Bias as a Specification Error", *Econometrica*, Vol. 47, pp. 153-161.
- HECKMAN, J. J. and R. ROBB, 1985. "Alternative Methods for Solving the Problem of Selection Bias in Evaluating the Impact of Treatments on Outcomes", in Wainer, H. (ed.), *Drawing Inferences from Self-Selected Samples*, Berlin: Springer-Verlag, pp. 63-107.
- HECKMAN, J. J. and T. E. MACURDY, 1986. "Labour Econometrics", in Griliches, Z. and M. D. Intriligator (eds.), *Handbook of Econometrics Vol. 3*, Amsterdam: North-Holland, pp. 1917-1977.
- HECKMAN, J. J. and V. J. HOLTZ, 1989. "Choosing among Alternative Nonexperimental Methods for Estimating the Impact of Social Programs: the Case of Manpower Training", *Journal of the*

- American Statistical Association*, Vol. 84, pp. 862-874.
- HECKMAN, J. J., 1990. "Varieties of Selection Bias", *American Economic Review*, Vol. 80, pp. 313-318.
- HECKMAN, J. J., 1996. "Comment on ANGRIST, J. D., G. W. IMBENS AND D.B. RUBIN", *Journal of the American Statistical Association*, Vol. 91, pp. 459-462.
- HECKMAN, J. J., H. ICHIMURA, J. SMITH and P. TODD, 1996. "Sources of Selection Bias in Evaluating Social Programs: an Interpretation of Conventional Measures and Evidence on the Effectiveness of Matching as a Program Evaluation Method", *Proceedings of the National Academy of Sciences USA*, Vol. 93, pp. 13416-13420.
- Using Experimental Data", *Econometrica*, Vol.66, pp. 1017-1098.
- HECKMAN, J. J. and J. A. SMITH, 1996. "Experimental and Non-experimental Evaluation", in Schmid, G, J. O Reilly and K. Schomann (eds.), *International Handbook of Labour Market Policy and Evaluation*, Cheltenham: Edward Elgar, pp. 37-88.
- HECKMAN, J. J., J. SMITH and N. CLEMENTS, 1997. "Making the Most of Programme Evaluations and Social Experiments: Accounting for Heterogeneity in Programme Impacts", *Review of Economic Studies*, Vol.64, pp. 487-535.
- HECKMAN, J. J., H. ICHIMURA and P. TODD, 1997. "Matching as an Econometric Evaluation Estimator: Evidence from Evaluating a Job Training Programme", *Review of Economic Studies*, Vol.64, pp. 605-654.
- HECKMAN, J. J., H. ICHIMURA and P. TODD, 1998. "Matching as an Econometric Evaluation Estimator", *Review of Economic Studies*, Vol.65, pp. 261-294.
- HECKMAN, J. J., H. ICHIMURA, J. SMITH and P. TODD, 1998. "Characterising Selection Bias Using Experimental Data", *Econometrica*, Vol.66, pp. 1017-1098.
- HOLLAND, P. W., 1989. "Comment on HECKMAN, J. J. and V. J. HOLTZ", *Journal of the American Statistical Association*, Vol. 84, pp. 875-877.
- HOROWITZ, J. L. and C. F. MANSKI, 1998. "Censoring of Outcomes and Regressors due to Survey Nonresponse: Identification and Estimation using Weights and Imputations", *Journal of Econometrics*, Vol. 84, pp. 37-58.
- IMBENS, G. W. and D.B. RUBIN, 1997. "Estimating Outcome Distributions for Compliers in Instrumental Variables Models", *Review of Economic Studies*, Vol. 64, pp. 555-574.
- JOHNSTON, J. and J. DINARDO, 1997. *Econometric Methods*, 4th Ed., New York: McGraw Hill.
- LALONDE, R., 1986. "Evaluating the Econometric Evaluations of Training Programmes with Experimental Data", *American Economic Review*, Vol. 76, pp. 604-620.
- LITTLE, R., 1985. "A Note about Models for Selectivity Bias", *Econometrica*, Vol. 53, pp. 1469-1474.
- LITTLE, R. and D. B. RUBIN, 1999. "Comment on Scharfstein, Rotnitzky and Robins", *Journal of the American Statistical Association*, Vol. 94, pp. 1130-1132.
- MADDALA, G. S., 1983. *Limited-Dependent and Qualitative Variables in Econometrics*, Cambridge: Cambridge University Press.
- MANSKI, C., 1990. "Nonparametric Bounds on Treatment Effects", *American Economic Review*, Vol. 80, pp. 319-323.
- NEWWEY, W. K., J. L. POWELL and J. WALKER, 1990. "Semi-parametric Estimation of Selection Models: Some Empirical Results", *American Economic Review*, Vol. 80, pp. 324-328.
- OBENCHAIN, R. L. and C. A. MELFI, 2000. "Propensity Score and Heckman Adjustments for Treatment Selection Bias in Database Studies", *American Statistical Association's 1999 Proceedings, Section on Statistics in Epidemiology*, to appear.
- O'CONNELL, P. and M. LYONS, 1995. *Enterprise-Related Training and State Policy in Ireland: The Training Support Scheme*, ESRI Policy Research Series Paper No. 25, Dublin: ESRI.
- O'CONNELL, P. and F. MCGINNITY, 1997. *Working Schemes? Active Labour Market Policy in Ireland*, Aldershot: Ashgate.
- PERKINS, S.M., W. TU, M. G. UNDERHILL, X-H ZHOU and M. D. MURRAY, 2000. "The Use of Propensity Scores in Pharmacoepidemiologic Research", *Pharmacoepidemiology and Drug Safety*, to appear.
- ROSENBAUM, P. R. and D. B. RUBIN, 1983a. "The Central Role of the Propensity Score in Observational Studies for Causal Effects", *Biometrika*, Vol. 70, pp. 41-55.
- ROSENBAUM, P. R. and D. B. RUBIN, 1983b. "Assessing Sensitivity to an Unobserved Binary Covariate in an Observational Study with Binary Outcome", *Journal of the Royal Statistical Society B*, Vol. 45, pp. 212-218.
- ROSENBAUM, P. R., 1984. "From Association to Causation in Observational Studies: the Role of Tests of Strongly Ignorable Treatment Assignment", *Journal of the American Statistical Association*, Vol. 79, pp. 41-47.

- ROSENBAUM, P. R. and D. B. RUBIN, 1984. "Reducing Bias in Observational Studies using Subclassification on the Propensity Score", *Journal of the American Statistical Association*, Vol. 79, pp. 516-524.
- ROSENBAUM, P. R. and D. B. RUBIN, 1985. "Constructing a Control Group Using a Multivariate Matched Sampling Method that Incorporates the Propensity Score", *The American Statistician*, Vol. 39, pp. 33-38.
- ROSENBAUM, P. R., 1987. "Model-Based Direct Adjustment", *Journal of the American Statistical Association*, Vol. 82, pp. 387-394.
- ROSENBAUM, P. R., 1989. "The Role of Known Effects in Observational Studies", *Biometrics*, Vol. 45, pp. 557-569.
- RUBIN, D. B., 1997. "Estimating Causal Effects from Large Data Sets using Propensity Scores", *Annals of Internal Medicine*, Vol. 127, pp. 757-763.
- RUBIN, D. B., H. STERN and V. VEHOVAR, 1995. "Handling 'Don't Know' Survey Responses; the Case of the Slovenian Plebiscite", *Journal of the American Statistical Association*, Vol. 90, pp. 822-828.
- TU, W., S. M. PERKINS, X-H ZHOU and M. D. MURRAY, 2000. "Testing Treatment Effect Using Propensity Score Stratification", Treatment Selection Bias in Database Studies", *American Statistical Association's 1999 Proceedings, Section on Statistics in Epidemiology*, to appear.