

# SOME PRACTICAL GUIDANCE FOR THE IMPLEMENTATION OF PROPENSITY SCORE MATCHING

Marco Caliendo

*IZA, Bonn*

Sabine Kopeinig

*University of Cologne*

**Abstract.** Propensity score matching (PSM) has become a popular approach to estimate causal treatment effects. It is widely applied when evaluating labour market policies, but empirical examples can be found in very diverse fields of study. Once the researcher has decided to use PSM, he is confronted with a lot of questions regarding its implementation. To begin with, a first decision has to be made concerning the estimation of the propensity score. Following that one has to decide which matching algorithm to choose and determine the region of common support. Subsequently, the matching quality has to be assessed and treatment effects and their standard errors have to be estimated. Furthermore, questions like ‘what to do if there is choice-based sampling?’ or ‘when to measure effects?’ can be important in empirical studies. Finally, one might also want to test the sensitivity of estimated treatment effects with respect to unobserved heterogeneity or failure of the common support condition. Each implementation step involves a lot of decisions and different approaches can be thought of. The aim of this paper is to discuss these implementation issues and give some guidance to researchers who want to use PSM for evaluation purposes.

**Keywords.** Propensity score matching; Treatment effects; Evaluation; Sensitivity analysis; Implementation

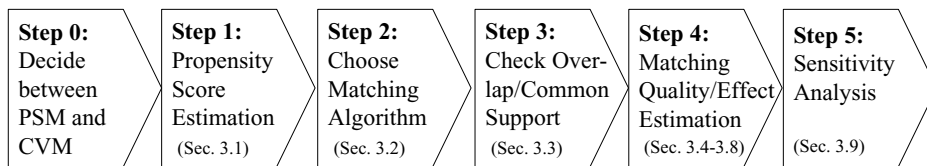
## 1. Introduction

Matching has become a popular approach to estimate causal treatment effects. It is widely applied when evaluating labour market policies (see e.g., Heckman *et al.*, 1997a; Dehejia and Wahba, 1999), but empirical examples can be found in very diverse fields of study. It applies for all situations where one has a treatment, a group of treated individuals and a group of untreated individuals. The nature of treatment may be very diverse. For example, Perkins *et al.* (2000) discuss the usage of matching in pharmacoepidemiologic research. Hitt and Frei (2002) analyse the effect of online banking on the profitability of customers. Davies and Kim (2003) compare

the effect on the percentage bid–ask spread of Canadian firms being interlisted on a US Exchange, whereas Brand and Halaby (2006) analyse the effect of elite college attendance on career outcomes. Ham *et al.* (2004) study the effect of a migration decision on the wage growth of young men and Bryson (2002) analyses the effect of union membership on wages of employees. Every microeconomic evaluation study has to overcome the fundamental evaluation problem and address the possible occurrence of selection bias. The first problem arises because we would like to know the difference between the participants' outcome with and without treatment. Clearly, we cannot observe both outcomes for the same individual at the same time. Taking the mean outcome of nonparticipants as an approximation is not advisable, since participants and nonparticipants usually differ even in the absence of treatment. This problem is known as selection bias and a good example is the case where high-skilled individuals have a higher probability of entering a training programme and also have a higher probability of finding a job. The matching approach is one possible solution to the selection problem. It originated from the statistical literature and shows a close link to the experimental context.<sup>1</sup> Its basic idea is to find in a large group of nonparticipants those individuals who are similar to the participants in all relevant pretreatment characteristics  $X$ . That being done, differences in outcomes of this well selected and thus adequate control group and of participants can be attributed to the programme. The underlying identifying assumption is known as unconfoundedness, selection on observables or conditional independence. It should be clear that matching is no 'magic bullet' that will solve the evaluation problem in any case. It should only be applied if the underlying identifying assumption can be credibly invoked based on the informational richness of the data and a detailed understanding of the institutional set-up by which selection into treatment takes place (see for example the discussion in Blundell *et al.*, 2005). For the rest of the paper we will assume that this assumption holds.

Since conditioning on all relevant covariates is limited in the case of a high dimensional vector  $X$  ('curse of dimensionality'), Rosenbaum and Rubin (1983b) suggest the use of so-called balancing scores  $b(X)$ , i.e. functions of the relevant observed covariates  $X$  such that the conditional distribution of  $X$  given  $b(X)$  is independent of assignment into treatment. One possible balancing score is the propensity score, i.e. the probability of participating in a programme given observed characteristics  $X$ . Matching procedures based on this balancing score are known as propensity score matching (PSM) and will be the focus of this paper. Once the researcher has decided to use PSM, he is confronted with a lot of questions regarding its implementation. Figure 1 summarizes the necessary steps when implementing PSM.<sup>2</sup>

The aim of this paper is to discuss these issues and give some practical guidance to researchers who want to use PSM for evaluation purposes. The paper is organized as follows. In Section 2, we will describe the basic evaluation framework and possible treatment effects of interest. Furthermore we show how PSM solves the evaluation problem and highlight the implicit identifying assumptions. In Section 3, we will focus on implementation steps of PSM estimators. To begin with, a first decision has to be made concerning the estimation of the propensity score



CVM: Covariate Matching, PSM: Propensity Score Matching

**Figure 1.** PSM – Implementation Steps.

(see Section 3.1). One has not only to decide about the probability model to be used for estimation, but also about variables which should be included in this model. In Section 3.2, we briefly evaluate the (dis-)advantages of different matching algorithms. Following that we discuss how to check the overlap between treatment and comparison group and how to implement the common support requirement in Section 3.3. In Section 3.4 we will show how to assess the matching quality. Subsequently we present the problem of choice-based sampling and discuss the question ‘when to measure programme effects?’ in Sections 3.5 and 3.6. Estimating standard errors for treatment effects will be discussed in Section 3.7 before we show in 3.8 how PSM can be combined with other evaluation methods. The following Section 3.9 is concerned with sensitivity issues, where we first describe approaches that allow researchers to determine the sensitivity of estimated effects with respect to a failure of the underlying unconfoundedness assumption. After that we introduce an approach that incorporates information from those individuals who failed the common support restriction, to calculate bounds of the parameter of interest, if all individuals from the sample at hand would have been included. Section 3.10 will briefly discuss the issues of programme heterogeneity, dynamic selection problems, and the choice of an appropriate control group and includes also a brief review of the available software to implement matching. Finally, Section 4 reviews all steps and concludes.

## 2. Evaluation Framework and Matching Basics

### *Roy–Rubin Model*

Inference about the impact of a treatment on the outcome of an individual involves speculation about how this individual would have performed had (s)he not received the treatment. The standard framework in evaluation analysis to formalize this problem is the potential outcome approach or Roy–Rubin model (Roy, 1951; Rubin, 1974). The main pillars of this model are individuals, treatment and potential outcomes. In the case of a binary treatment the treatment indicator  $D_i$  equals one if individual  $i$  receives treatment and zero otherwise. The potential outcomes are then defined as  $Y_i(D_i)$  for each individual  $i$ , where  $i = 1, \dots, N$  and  $N$  denotes the total population. The treatment effect for an individual  $i$  can be written as

$$\tau_i = Y_i(1) - Y_i(0) \quad (1)$$

The fundamental evaluation problem arises because only one of the potential outcomes is observed for each individual  $i$ . The unobserved outcome is called the counterfactual outcome. Hence, estimating the individual treatment effect  $\tau_i$  is not possible and one has to concentrate on (population) average treatment effects.<sup>3</sup>

### *Parameter of Interest and Selection Bias*

Two parameters are most frequently estimated in the literature. The first one is the population average treatment effect (ATE), which is simply the difference of the expected outcomes after participation and nonparticipation:

$$\tau_{ATE} = E(\tau) = E[Y(1) - Y(0)] \quad (2)$$

This parameter answers the question: ‘What is the expected effect on the outcome if individuals in the population were randomly assigned to treatment?’ Heckman (1997) notes that this estimate might not be of relevance to policy makers because it includes the effect on persons for whom the programme was never intended. For example, if a programme is specifically targeted at individuals with low family income, there is little interest in the effect of such a programme for a millionaire. Therefore, the most prominent evaluation parameter is the so-called average treatment effect on the treated (ATT), which focuses explicitly on the effects on those for whom the programme is actually intended. It is given by

$$\tau_{ATT} = E(\tau|D = 1) = E[Y(1)|D = 1] - E[Y(0)|D = 1] \quad (3)$$

The expected value of ATT is defined as the difference between expected outcome values with and without treatment for those who actually participated in treatment. In the sense that this parameter focuses directly on actual treatment participants, it determines the realized gross gain from the programme and can be compared with its costs, helping to decide whether the programme is successful or not (Heckman *et al.*, 1999). The most interesting parameter to estimate depends on the specific evaluation context and the specific question asked. Heckman *et al.* (1999) discuss further parameters, like the proportion of participants who benefit from the programme or the distribution of gains at selected base state values. For most evaluation studies, however, the focus lies on ATT and therefore we will focus on this parameter, too.<sup>4</sup> As the counterfactual mean for those being treated –  $E[Y(0)|D = 1]$  – is not observed, one has to choose a proper substitute for it in order to estimate ATT. Using the mean outcome of untreated individuals  $E[Y(0)|D = 0]$  is in nonexperimental studies usually not a good idea, because it is most likely that components which determine the treatment decision also determine the outcome variable of interest. Thus, the outcomes of individuals from the treatment and comparison groups would differ even in the absence of treatment leading to a ‘selection bias’. For ATT it can be noted as

$$E[Y(1)|D = 1] - E[Y(0)|D = 0] = \tau_{ATT} + E[Y(0)|D = 1] - E[Y(0)|D = 0] \quad (4)$$

The difference between the left-hand side of equation (4) and  $\tau_{ATT}$  is the so-called ‘selection bias’. The true parameter  $\tau_{ATT}$  is only identified if

$$E[Y(0)|D = 1] - E[Y(0)|D = 0] = 0 \quad (5)$$

In social experiments where assignment to treatment is random this is ensured and the treatment effect is identified.<sup>5</sup> In nonexperimental studies one has to invoke some identifying assumptions to solve the selection problem stated in equation (4).

### *Unconfoundedness and Common Support*

One major strand of evaluation literature focuses on the estimation of treatment effects under the assumption that the treatment satisfies some form of exogeneity. Different versions of this assumption are referred to as unconfoundedness (Rosenbaum and Rubin, 1983b), selection on observables (Heckman and Robb, 1985) or conditional independence assumption (CIA) (Lechner, 1999). We will use these terms throughout the paper interchangeably. This assumption implies that systematic differences in outcomes between treated and comparison individuals with the same values for covariates are attributable to treatment. Imbens (2004) gives an extensive overview of estimating ATEs under unconfoundedness. The identifying assumption can be written as

**Assumption 1.** *Unconfoundedness:*  $Y(0), Y(1) \perp\!\!\!\perp D \mid X$

where  $\perp\!\!\!\perp$  denotes independence, i.e. given a set of observable covariates  $X$  which are not affected by treatment, potential outcomes are independent of treatment assignment. This implies that all variables that influence treatment assignment and potential outcomes simultaneously have to be observed by the researcher. Clearly, this is a strong assumption and has to be justified by the data quality at hand. For the rest of the paper we will assume that this condition holds. If the researcher believes that the available data are not rich enough to justify this assumption, he has to rely on different identification strategies which explicitly allow selection on unobservables, too. Prominent examples are difference-in-differences (DID) and instrumental variables estimators.<sup>6</sup> We will show in Section 3.8 how propensity score matching can be combined with some of these methods.

A further requirement besides independence is the common support or overlap condition. It rules out the phenomenon of perfect predictability of  $D$  given  $X$ .

**Assumption 2.** *Overlap:*  $0 < P(D = 1|X) < 1$ .

It ensures that persons with the same  $X$  values have a positive probability of being both participants and nonparticipants (Heckman *et al.*, 1999). Rosenbaum and Rubin (1983b) call Assumptions 1 and 2 together ‘strong ignorability’. Under ‘strong ignorability’ ATE in (2) and ATT in (3) can be defined for all values of  $X$ . Heckman *et al.* (1998b) demonstrate that the ignorability or unconfoundedness conditions are overly strong. All that is needed for estimation of (2) and (3) is mean independence. However, Lechner (2002) argues that Assumption 1 has the virtue of identifying

mean effects for all transformations of the outcome variables. The reason is that the weaker assumption of mean independence is intrinsically tied to functional form assumptions, making an identification of average effects on transformations of the original outcome impossible (Imbens, 2004). Furthermore, it will be difficult to argue why conditional mean independence should hold and Assumption 1 might still be violated in empirical studies.

If we are interested in estimating the ATT only, we can weaken the unconfoundedness assumption in a different direction. In that case one needs only to assume

**Assumption 3.** *Unconfoundedness for controls:*  $Y(0) \perp\!\!\!\perp D \mid X$

and the weaker overlap assumption

**Assumption 4.** *Weak overlap:*  $P(D = 1 \mid X) < 1$ .

These assumptions are sufficient for identification of (3), because the moments of the distribution of  $Y(1)$  for the treated are directly estimable.

#### *Unconfoundedness given the Propensity Score*

It should also be clear that conditioning on all relevant covariates is limited in the case of a high dimensional vector  $X$ . For instance if  $X$  contains  $s$  covariates which are all dichotomous, the number of possible matches will be  $2^s$ . To deal with this dimensionality problem, Rosenbaum and Rubin (1983b) suggest using so-called balancing scores. They show that if potential outcomes are independent of treatment conditional on covariates  $X$ , they are also independent of treatment conditional on a balancing score  $b(X)$ . The propensity score  $P(D = 1 \mid X) = P(X)$ , i.e. the probability for an individual to participate in a treatment given his observed covariates  $X$ , is one possible balancing score. Hence, if Assumption 1 holds, all biases due to observable components can be removed by conditioning on the propensity score (Imbens, 2004).

**Corollary 1.** *Unconfoundedness given the propensity score:*  $Y(0), Y(1) \perp\!\!\!\perp D \mid P(X)$ .<sup>7</sup>

#### *Estimation Strategy*

Given that CIA holds and assuming additionally that there is overlap between both groups, the PSM estimator for ATT can be written in general as

$$\tau_{ATT}^{PSM} = E_{P(X) \mid D=1} \{E[Y(1) \mid D = 1, P(X)] - E[Y(0) \mid D = 0, P(X)]\} \quad (6)$$

To put it in words, the PSM estimator is simply the mean difference in outcomes over the common support, appropriately weighted by the propensity score distribution of participants. Based on this brief outline of the matching estimator in the general evaluation framework, we are now going to discuss the implementation of PSM in detail.

### 3. Implementation of Propensity Score Matching

#### 3.1 *Estimating the Propensity Score*

When estimating the propensity score, two choices have to be made. The first one concerns the model to be used for the estimation, and the second one the variables to be included in this model. We will start with the model choice before we discuss which variables to include in the model.

##### *Model Choice – Binary Treatment*

Little advice is available regarding which functional form to use (see for example the discussion in Smith, 1997). In principle any discrete choice model can be used. Preference for logit or probit models (compared to linear probability models) derives from the well-known shortcomings of the linear probability model, especially the unlikeliness of the functional form when the response variable is highly skewed and predictions that are outside the  $[0, 1]$  bounds of probabilities. However, when the purpose of a model is classification rather than estimation of structural coefficients, it is less clear that these criticisms apply (Smith, 1997). For the binary treatment case, where we estimate the probability of participation versus nonparticipation, logit and probit models usually yield similar results. Hence, the choice is not too critical, even though the logit distribution has more density mass in the bounds.

##### *Model Choice – Multiple Treatments*

However, when leaving the binary treatment case, the choice of the model becomes more important. The multiple treatment case (as discussed in Imbens (2000) and Lechner (2001a)) consists of more than two alternatives, for example when an individual is faced with the choice to participate in job-creation schemes, vocational training or wage subsidy programmes or to not participate at all (we will describe this approach in more detail in Section 3.10). For that case it is well known that the multinomial logit is based on stronger assumptions than the multinomial probit model, making the latter the preferable option.<sup>8</sup> However, since the multinomial probit is computationally more burdensome, a practical alternative is to estimate a series of binomial models as suggested by Lechner (2001a). Bryson *et al.* (2002) note that there are two shortcomings regarding this approach. First, as the number of options increases, the number of models to be estimated increases disproportionately (for  $L$  options we need  $0.5(L(L - 1))$  models). Second, in each model only two options at a time are considered and consequently the choice is conditional on being in one of the two selected groups. On the other hand, Lechner (2001a) compares the performance of the multinomial probit approach and series estimation and finds little difference in their relative performance. He suggests that the latter approach may be more robust since a mis-specification in one of the series will not compromise all others as would be the case in the multinomial probit model.

### *Variable Choice:*

More advice is available regarding the inclusion (or exclusion) of covariates in the propensity score model. The matching strategy builds on the CIA, requiring that the outcome variable(s) must be independent of treatment conditional on the propensity score. Hence, implementing matching requires choosing a set of variables  $X$  that credibly satisfy this condition. Heckman *et al.* (1997a) and Dehejia and Wahba (1999) show that omitting important variables can seriously increase bias in resulting estimates. Only variables that influence simultaneously the participation decision and the outcome variable should be included. Hence, economic theory, a sound knowledge of previous research and also information about the institutional settings should guide the researcher in building up the model (see e.g., Sianesi, 2004; Smith and Todd, 2005). It should also be clear that only variables that are unaffected by participation (or the anticipation of it) should be included in the model. To ensure this, variables should either be fixed over time or measured before participation. In the latter case, it must be guaranteed that the variable has not been influenced by the anticipation of participation. Heckman *et al.* (1999) also point out that the data for participants and nonparticipants should stem from the same sources (e.g. the same questionnaire). The better and more informative the data are, the easier it is to credibly justify the CIA and the matching procedure. However, it should also be clear that ‘too good’ data is not helpful either. If  $P(X) = 0$  or  $P(X) = 1$  for some values of  $X$ , then we cannot use matching conditional on those  $X$  values to estimate a treatment effect, because persons with such characteristics either always or never receive treatment. Hence, the common support condition as stated in Assumption 2 fails and matches cannot be performed. Some randomness is needed that guarantees that persons with identical characteristics can be observed in both states (Heckman *et al.*, 1998b).

In cases of uncertainty of the proper specification, sometimes the question may arise whether it is better to include too many rather than too few variables. Bryson *et al.* (2002) note that there are two reasons why over-parameterized models should be avoided. First, it may be the case that including extraneous variables in the participation model exacerbates the support problem. Second, although the inclusion of nonsignificant variables in the propensity score specification will not bias the propensity score estimates or make them inconsistent, it can increase their variance.

The results from Augurzky and Schmidt (2001) point in the same direction. They run a simulation study to investigate PSM when selection into treatment is remarkably strong, and treated and untreated individuals differ considerably in their observable characteristics. In their set-up, explanatory variables in the selection equation are partitioned into three sets. The first set (set 1) includes covariates which strongly influence the treatment decision but weakly influence the outcome variable. Furthermore, they include covariates which are relevant to the outcome but irrelevant to the treatment decision (set 2) and covariates which influence both (set 3). Including the full set of covariates in the propensity score specification (full model including all three sets of covariates) might cause problems in small samples in terms of higher variance, since either some treated have to be discarded from the



analysis or control units have to be used more than once. They show that matching on an inconsistent estimate of the propensity score (i.e. partial model including only set 3 or both sets 1 and 3) produces better estimation results of the ATE.

On the other hand, Rubin and Thomas (1996) recommend against ‘trimming’ models in the name of parsimony. They argue that a variable should only be excluded from analysis if there is consensus that the variable is either unrelated to the outcome or not a proper covariate. If there are doubts about these two points, they explicitly advise to include the relevant variables in the propensity score estimation.

By these criteria, there are both reasons for and against including all of the reasonable covariates available. Basically, the points made so far imply that the choice of variables should be based on economic theory and previous empirical findings. But clearly, there are also some formal (statistical) tests which can be used. Heckman *et al.* (1998a), Heckman and Smith (1999) and Black and Smith (2004) discuss three strategies for the selection of variables to be used in estimating the propensity score.

### *Hit or Miss Method*

The first one is the ‘hit or miss’ method or prediction rate metric, where variables are chosen to maximize the within-sample correct prediction rates. This method classifies an observation as ‘1’ if the estimated propensity score is larger than the sample proportion of persons taking treatment, i.e.  $\hat{P}(X) > \bar{P}$ . If  $\hat{P}(X) \leq \bar{P}$  observations are classified as ‘0’. This method maximizes the overall classification rate for the sample assuming that the costs for the misclassification are equal for the two groups (Heckman *et al.*, 1997a).<sup>9</sup> But clearly, it has to be kept in mind that the main purpose of the propensity score estimation is not to predict selection into treatment as well as possible but to balance all covariates (Augurzy and Schmidt, 2001).

### *Statistical Significance*

The second approach relies on statistical significance and is very common in textbook econometrics. To do so, one starts with a parsimonious specification of the model, e.g. a constant, the age and some regional information, and then ‘tests up’ by iteratively adding variables to the specification. A new variable is kept if it is statistically significant at conventional levels. If combined with the ‘hit or miss’ method, variables are kept if they are statistically significant and increase the prediction rates by a substantial amount (Heckman *et al.*, 1998a).

### *Leave-One-Out Cross-Validation*

Leave-one-out cross-validation can also be used to choose the set of variables to be included in the propensity score. Black and Smith (2004) implement their model selection procedure by starting with a ‘minimal’ model containing only two variables. They subsequently add blocks of additional variables and compare the

resulting mean squared errors. As a note of caution, they stress that this amounts to choosing the propensity score model based on goodness-of-fit considerations, rather than based on theory and evidence about the set of variables related to the participation decision and the outcomes (Black and Smith, 2004). They also point out an interesting trade-off in finite samples between the plausibility of the CIA and the variance of the estimates. When using the full specification, bias arises from selecting a wide bandwidth in response to the weakness of the common support. In contrast to that, when matching on the minimal specification, common support is not a problem but the plausibility of the CIA is. This trade-off also affects the estimated standard errors, which are smaller for the minimal specification where the common support condition poses no problem.

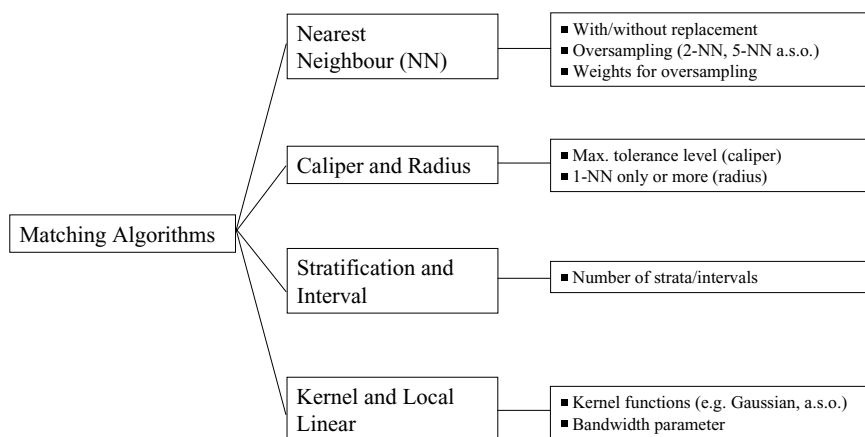
Finally, checking the matching quality can also help to determine the propensity score specification and we will discuss this point later in Section 3.4.

### *Overweighting some Variables*

Let us assume for the moment that we have found a satisfactory specification of the model. It may sometimes be felt that some variables play a specifically important role in determining participation and outcome (Bryson *et al.*, 2002). As an example, one can think of the influence of gender and region in determining the wage of individuals. Let us take as given for the moment that men earn more than women and the wage level is higher in region A compared to region B. If we add dummy variables for gender and region in the propensity score estimation, it is still possible that women in region B are matched with men in region A, since the gender and region dummies are only a subset of all available variables. There are basically two ways to put greater emphasis on specific variables. One can either find variables in the comparison group who are identical with respect to these variables, or carry out matching on subpopulations. The study from Lechner (2002) is a good example for the first approach. He evaluates the effects of active labour market policies in Switzerland and uses the propensity score as a 'partial' balancing score which is complemented by an exact matching on sex, duration of unemployment and native language. Heckman *et al.* (1997a, 1998a) use the second strategy and implement matching separately for four demographic groups. That implies that the complete matching procedure (estimating the propensity score, checking the common support, etc.) has to be implemented separately for each group. This is analogous to insisting on a perfect match, e.g. in terms of gender and region, and then carrying out propensity score matching. This procedure is especially recommendable if one expects the effects to be heterogeneous between certain groups.

### *Alternatives to the Propensity Score*

Finally, it should also be noted that it is possible to match on a measure other than the propensity score, namely the underlying index of the score estimation. The advantage of this is that the index differentiates more between observations in the extremes of the distribution of the propensity score (Lechner, 2000). This is



**Figure 2.** Different Matching Algorithms.

useful if there is some concentration of observations in the tails of the distribution. Additionally, in some recent papers the propensity score is estimated by duration models. This is of particular interest if the ‘timing of events’ plays a crucial role (see e.g. Brodaty *et al.*, 2001; Sianesi, 2004).

### 3.2 Choosing a Matching Algorithm

The PSM estimator in its general form was stated in equation (6). All matching estimators contrast the outcome of a treated individual with outcomes of comparison group members. PSM estimators differ not only in the way the neighbourhood for each treated individual is defined and the common support problem is handled, but also with respect to the weights assigned to these neighbours. Figure 2 depicts different PSM estimators and the inherent choices to be made when they are used. We will not discuss the technical details of each estimator here at depth but rather present the general ideas and the involved trade-offs with each algorithm.<sup>10</sup>

#### *Nearest Neighbour Matching*

The most straightforward matching estimator is nearest neighbour (NN) matching. The individual from the comparison group is chosen as a matching partner for a treated individual that is closest in terms of the propensity score. Several variants of NN matching are proposed, e.g. NN matching ‘with replacement’ and ‘without replacement’. In the former case, an untreated individual can be used more than once as a match, whereas in the latter case it is considered only once. Matching with replacement involves a trade-off between bias and variance. If we allow replacement, the average quality of matching will increase and the bias will decrease. This is of

particular interest with data where the propensity score distribution is very different in the treatment and the control group. For example, if we have a lot of treated individuals with high propensity scores but only few comparison individuals with high propensity scores, we get bad matches as some of the high-score participants will get matched to low-score nonparticipants. This can be overcome by allowing replacement, which in turn reduces the number of distinct nonparticipants used to construct the counterfactual outcome and thereby increases the variance of the estimator (Smith and Todd, 2005). A problem which is related to NN matching without replacement is that estimates depend on the order in which observations get matched. Hence, when using this approach it should be ensured that ordering is randomly done.<sup>11</sup>

It is also suggested to use more than one NN ('oversampling'). This form of matching involves a trade-off between variance and bias, too. It trades reduced variance, resulting from using more information to construct the counterfactual for each participant, with increased bias that results from on average poorer matches (see e.g. Smith, 1997). When using oversampling, one has to decide how many matching partners should be chosen for each treated individual and which weight (e.g. uniform or triangular weight) should be assigned to them.

### *Caliper and Radius Matching*

NN matching faces the risk of bad matches if the closest neighbour is far away. This can be avoided by imposing a tolerance level on the maximum propensity score distance (caliper). Hence, caliper matching is one form of imposing a common support condition (we will come back to this point in Section 3.3). Bad matches are avoided and the matching quality rises. However, if fewer matches can be performed, the variance of the estimates increases.<sup>12</sup> Applying caliper matching means that an individual from the comparison group is chosen as a matching partner for a treated individual that lies within the caliper ('propensity range') and is closest in terms of propensity score. As Smith and Todd (2005) note, a possible drawback of caliper matching is that it is difficult to know *a priori* what choice for the tolerance level is reasonable.

Dehejia and Wahba (2002) suggest a variant of caliper matching which is called radius matching. The basic idea of this variant is to use not only the NN within each caliper but all of the comparison members within the caliper. A benefit of this approach is that it uses only as many comparison units as are available within the caliper and therefore allows for usage of extra (fewer) units when good matches are (not) available. Hence, it shares the attractive feature of oversampling mentioned above, but avoids the risk of bad matches.

### *Stratification and Interval Matching*

The idea of stratification matching is to partition the common support of the propensity score into a set of intervals (strata) and to calculate the impact within each interval by taking the mean difference in outcomes between treated and

control observations. This method is also known as interval matching, blocking and subclassification (Rosenbaum and Rubin, 1984). Clearly, one question to be answered is how many strata should be used in empirical analysis. Cochran (1968) shows that five subclasses are often enough to remove 95% of the bias associated with one single covariate. Since, as Imbens (2004) notes, all bias under unconfoundedness is associated with the propensity score, this suggests that under normality the use of five strata removes most of the bias associated with all covariates. One way to justify the choice of the number of strata is to check the balance of the propensity score (or the covariates) within each stratum (see e.g. Aakvik, 2001). Most of the algorithms can be described in the following way. First, check if within a stratum the propensity score is balanced. If not, strata are too large and need to be split. If, conditional on the propensity score being balanced, the covariates are unbalanced, the specification of the propensity score is not adequate and has to be respecified, e.g. through the addition of higher-order terms or interactions (see Dehejia and Wahba, 1999; Dehejia, 2005).

### *Kernel and Local Linear Matching*

The matching algorithms discussed so far have in common that only a few observations from the comparison group are used to construct the counterfactual outcome of a treated individual. Kernel matching (KM) and local linear matching (LLM) are nonparametric matching estimators that use weighted averages of (nearly) all – depending on the choice of the kernel function – individuals in the control group to construct the counterfactual outcome. Thus, one major advantage of these approaches is the lower variance which is achieved because more information is used. A drawback of these methods is that possibly observations are used that are bad matches. Hence, the proper imposition of the common support condition is of major importance for KM and LLM. Heckman *et al.* (1998b) derive the asymptotic distribution of these estimators and Heckman *et al.* (1997a) present an application. As Smith and Todd (2005) note, KM can be seen as a weighted regression of the counterfactual outcome on an intercept with weights given by the kernel weights. Weights depend on the distance between each individual from the control group and the participant observation for which the counterfactual is estimated. It is worth noting that if weights from a symmetric, nonnegative, unimodal kernel are used, then the average places higher weight on persons close in terms of the propensity score of a treated individual and lower weight on more distant observations. The estimated intercept provides an estimate of the counterfactual mean. The difference between KM and LLM is that the latter includes in addition to the intercept a linear term in the propensity score of a treated individual. This is an advantage whenever comparison group observations are distributed asymmetrically around the treated observation, e.g. at boundary points, or when there are gaps in the propensity score distribution. When applying KM one has to choose the kernel function and the bandwidth parameter. The first point appears to be relatively unimportant in practice (DiNardo and Tobias, 2001). What is seen as more important (see e.g. Silverman, 1986; Pagan and Ullah, 1999) is the choice of the bandwidth parameter with the following

**Table 1.** Trade-offs in Terms of Bias and Efficiency.

Decision	Bias	Variance
Nearest neighbour matching:		
multiple neighbours/single neighbour	(+)/(−)	(−)/(+)
with caliper/without caliper	(−)/(+)	(+)/(−)
Use of control individuals:		
with replacement/without replacement	(−)/(+)	(+)/(−)
Choosing method:		
NN matching/Radius matching	(−)/(+)	(+)/(−)
KM or LLM/NN methods	(+)/(−)	(−)/(+)
Bandwidth choice with KM:		
small/large	(−)/(+)	(+)/(−)
Polynomial order with LPM:		
small/large	(+)/(−)	(−)/(+)

KM, kernel matching, LLM; local linear matching; LPM, local polynomial matching NN, nearest neighbour; increase; (+); decrease (−).

trade-off arising. High bandwidth values yield a smoother estimated density function, therefore leading to a better fit and a decreasing variance between the estimated and the true underlying density function. On the other hand, underlying features may be smoothed away by a large bandwidth leading to a biased estimate. The bandwidth choice is therefore a compromise between a small variance and an unbiased estimate of the true density function. It should be noted that LLM is a special case of local polynomial matching (LPM). LPM includes in addition to an intercept a term of polynomial order  $p$  in the propensity score, e.g.  $p = 1$  for LLM,  $p = 2$  for local quadratic matching or  $p = 3$  for local cubic matching. Generally, the larger the polynomial order  $p$  is the smaller will be the asymptotic bias but the larger will be the asymptotic variance. To our knowledge Ham *et al.* (2004) is the only application of local cubic matching so far, and hence practical experiences with LPM estimators with  $p \geq 2$  are rather limited.

### *Trade-offs in Terms of Bias and Efficiency*

Having presented the different possibilities, the question remains of how one should select a specific matching algorithm. Clearly, asymptotically all PSM estimators should yield the same results, because with growing sample size they all become closer to comparing only exact matches (Smith, 2000). However, in small samples the choice of the matching algorithm can be important (Heckman *et al.*, 1997a), where usually a trade-off between bias and variance arises (see Table 1). So what advice can be given to researchers facing the problem of choosing a matching estimator? It

should be clear that there is no ‘winner’ for all situations and that the choice of the estimator crucially depends on the situation at hand. The performance of different matching estimators varies case-by-case and depends largely on the data structure at hand (Zhao, 2000). To give an example, if there are only a few control observations, it makes no sense to match without replacement. On the other hand, if there are a lot of comparable untreated individuals it might be worth using more than one NN (either by oversampling or KM) to gain more precision in estimates. Pragmatically, it seems sensible to try a number of approaches. Should they give similar results, the choice may be unimportant. Should results differ, further investigation may be needed in order to reveal more about the source of the disparity (Bryson *et al.*, 2002).

### 3.3 *Overlap and Common Support*

Our discussion in Section 2 has shown that ATT and ATE are only defined in the region of common support. Heckman *et al.* (1997a) point out that a violation of the common support condition is a major source of evaluation bias as conventionally measured. Comparing the incomparable must be avoided, i.e. only the subset of the comparison group that is comparable to the treatment group should be used in the analysis (Dehejia and Wahba, 1999). Hence, an important step is to check the overlap and the region of common support between treatment and comparison group. Several ways are suggested in the literature, where the most straightforward one is a visual analysis of the density distribution of the propensity score in both groups. Lechner (2001b) argues that given that the support problem can be spotted by inspecting the propensity score distribution, there is no need to implement a complicated estimator. However, some guidelines might help the researcher to determine the region of common support more precisely. We will present two methods, where the first one is essentially based on comparing the minima and maxima of the propensity score in both groups and the second one is based on estimating the density distribution in both groups. Implementing the common support condition ensures that any combination of characteristics observed in the treatment group can also be observed among the control group (Bryson *et al.*, 2002). For ATT it is sufficient to ensure the existence of potential matches in the control group, whereas for ATE it is additionally required that the combinations of characteristics in the comparison group may also be observed in the treatment group (Bryson *et al.*, 2002).

#### *Minima and Maxima Comparison*

The basic criterion of this approach is to delete all observations whose propensity score is smaller than the minimum and larger than the maximum in the opposite group. To give an example let us assume for a moment that the propensity score lies within the interval [0.07, 0.94] in the treatment group and within [0.04, 0.89] in the control group. Hence, with the ‘minima and maxima criterion’, the common support is given by [0.07, 0.89]. Observations which lie outside this region are discarded from analysis. Clearly a two-sided test is only necessary if the parameter

of interest is ATE; for ATT it is sufficient to ensure that for each participant a close nonparticipant can be found. It should also be clear that the common support condition is in some ways more important for the implementation of KM than it is for the implementation of NN matching, because with KM all untreated observations are used to estimate the missing counterfactual outcome, whereas with NN matching only the closest neighbour is used. Hence, NN matching (with the additional imposition of a maximum allowed caliper) handles the common support problem pretty well. There are some problems associated with the ‘minima and maxima comparison’, e.g. if there are observations at the bounds which are discarded even though they are very close to the bounds. Another problem arises if there are areas within the common support interval where there is only limited overlap between both groups, e.g. if in the region  $[0.51, 0.55]$  only treated observations can be found. Additionally problems arise if the density in the tails of the distribution is very thin, for example when there is a substantial distance from the smallest maximum to the second smallest element. Therefore, Lechner (2002) suggests to check the sensitivity of the results when the minima and maxima are replaced by the tenth smallest and tenth largest observation.

### *Trimming to Determine the Common Support*

A different way to overcome these possible problems is described by Smith and Todd (2005).<sup>13</sup> They use a trimming procedure to determine the common support region and define the region of common support as those values of  $P$  that have positive density within both the  $D = 1$  and  $D = 0$  distributions, i.e.

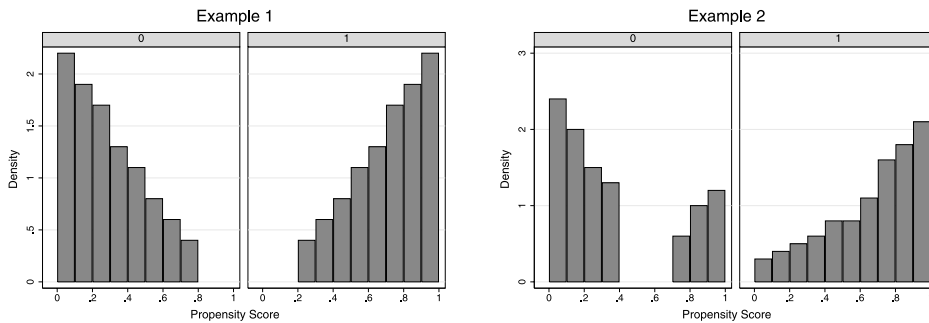
$$\hat{S}_P = \{P : \hat{f}(P|D = 1) > 0 \text{ and } \hat{f}(P|D = 0) > 0\} \quad (7)$$

where  $\hat{f}(P|D = 1) > 0$  and  $\hat{f}(P|D = 0) > 0$  are nonparametric density estimators. Any  $P$  points for which the estimated density is exactly zero are excluded. Additionally – to ensure that the densities are strictly positive – they require that the densities exceed zero by a threshold amount  $q$ . So not only the  $P$  points for which the estimated density is exactly zero, but also an additional  $q$  percent of the remaining  $P$  points for which the estimated density is positive but very low are excluded.<sup>14</sup>

$$\hat{S}_{Pq} = \{Pq : \hat{f}(P|D = 1) > q \text{ and } \hat{f}(P|D = 0) > q\} \quad (8)$$

Figure 3 gives a hypothetical example and clarifies the differences between the two approaches. In the first example the propensity score distribution is highly skewed to the left (right) for participants (nonparticipants). Even though this is an extreme example, researchers are confronted with similar distributions in practice, too. With the ‘minima and maxima comparison’ we would exclude any observations lying outside the region of common support given by  $[0.2, 0.8]$ . Depending on the chosen trimming level  $q$ , we would maybe also exclude control observations in the interval  $[0.7, 0.8]$  and treated observations in the interval  $[0.2, 0.3]$  with the trimming approach since the densities are relatively low there. However, no large differences between the two approaches would emerge. In the second example we





The left side in each example refers to non-participants ( $D=0$ ), the right side to participants ( $D=1$ ).

Source: Hypothetical example

**Figure 3.** The Common Support Problem.

do not find any control individuals in the region  $[0.4, 0.7]$ . The ‘minima and maxima comparison’ fails in that situation, since minima and maxima in both groups are equal at 0.01 and 0.99. Hence, no observations would be excluded based on this criterion making the estimation of treatment effects in the region  $[0.4, 0.7]$  questionable. The trimming method on the other hand would explicitly exclude treated observations in that propensity score range and would therefore deliver more reliable results.<sup>15</sup> Hence, the choice of the method depends on the data situation at hand and before making any decisions a visual analysis is recommended.

### *Failure of the Common Support*

Once one has defined the region of common support, individuals that fall outside this region have to be disregarded and for these individuals the treatment effect cannot be estimated. Bryson *et al.* (2002) note that when the proportion of lost individuals is small, this poses few problems. However, if the number is too large, there may be concerns whether the estimated effect on the remaining individuals can be viewed as representative. It may be instructive to inspect the characteristics of discarded individuals since those can provide important clues when interpreting the estimated treatment effects. Lechner (2001b) notes that both ignoring the support problem and estimating treatment effects only within the common support (subgroup effects) may be misleading. He develops an approach that can be used to derive bounds for the true treatment effect and we describe this approach in detail in Section 3.9.

### *3.4 Assessing the Matching Quality*

Since we do not condition on all covariates but on the propensity score, it has to be checked if the matching procedure is able to balance the distribution of the relevant variables in both the control and treatment group. Several procedures to do so will be discussed in this section. These procedures can also, as already mentioned, help in

determining which interactions and higher-order terms to include in the propensity score specification for a given set of covariates  $X$ . The basic idea of all approaches is to compare the situation before and after matching and check if there remain any differences after conditioning on the propensity score. If there are differences, matching on the score was not (completely) successful and remedial measures have to be done, e.g. by including interaction terms in the estimation of the propensity score. A helpful theorem in this context is suggested by Rosenbaum and Rubin (1983b) and states that

$$X \perp\!\!\!\perp D | P(D = 1|X) \quad (9)$$

This means that after conditioning on  $P(D = 1|X)$ , additional conditioning on  $X$  should not provide new information about the treatment decision. Hence, if after conditioning on the propensity score there is still dependence on  $X$ , this suggests either mis-specification in the model used to estimate  $P(D = 1|X)$  (see Smith and Todd, 2005) or a fundamental lack of comparability between the two groups (Blundell *et al.*, 2005).<sup>16</sup>

### Standardized Bias

One suitable indicator to assess the distance in marginal distributions of the  $X$  variables is the standardized bias (SB) suggested by Rosenbaum and Rubin (1985). For each covariate  $X$  it is defined as the difference of sample means in the treated and matched control subsamples as a percentage of the square root of the average of sample variances in both groups. The SB before matching is given by

$$SB_{\text{before}} = 100 \cdot \frac{\bar{X}_1 - \bar{X}_0}{\sqrt{0.5 \cdot (V_1(X) + V_0(X))}} \quad (10)$$

The SB after matching is given by

$$SB_{\text{after}} = 100 \cdot \frac{\bar{X}_{1M} - \bar{X}_{0M}}{\sqrt{0.5 \cdot (V_{1M}(X) + V_{0M}(X))}} \quad (11)$$

where  $X_1$  ( $V_1$ ) is the mean (variance) in the treatment group before matching and  $X_0$  ( $V_0$ ) the analogue for the control group.  $X_{1M}$  ( $V_{1M}$ ) and  $X_{0M}$  ( $V_{0M}$ ) are the corresponding values for the matched samples. This is a common approach used in many evaluation studies, e.g. by Lechner (1999), Sianesi (2004) and Caliendo *et al.* (2007). One possible problem with the SB approach is that one does not have a clear indication for the success of the matching procedure, even though in most empirical studies an SB below 3% or 5% after matching is seen as sufficient.

### *t*-Test

A similar approach uses a two-sample *t*-test to check if there are significant differences in covariate means for both groups (Rosenbaum and Rubin, 1985). Before matching differences are expected, but after matching the covariates should be balanced in both groups and hence no significant differences should be found. The

$t$ -test might be preferred if the evaluator is concerned with the statistical significance of the results. The shortcoming here is that the bias reduction before and after matching is not clearly visible.

### *Joint Significance and Pseudo- $R^2$*

Additionally, Sianesi (2004) suggests to reestimate the propensity score on the matched sample, i.e. only on participants and matched nonparticipants, and compare the pseudo- $R^2$ s before and after matching. The pseudo- $R^2$  indicates how well the regressors  $X$  explain the participation probability. After matching there should be no systematic differences in the distribution of covariates between both groups and therefore the pseudo- $R^2$  should be fairly low. Furthermore, one can also perform a likelihood ratio test on the joint significance of all regressors in the probit or logit model. The test should not be rejected before, and should be rejected after, matching.

### *Stratification Test*

Finally, Dehejia and Wahba (1999, 2002) divide observations into strata based on the estimated propensity score, such that no statistically significant difference between the mean of the estimated propensity score in both treatment and control group remain. Then they use  $t$ -tests within each strata to test if the distribution of  $X$  variables is the same between both groups (for the first and second moments). If there are remaining differences, they add higher-order and interaction terms in the propensity score specification, until such differences no longer emerge.

This makes clear that an assessment of matching quality can also be used to determine the propensity score specification. If the quality indicators are not satisfactory, one reason might be mis-specification of the propensity score model and hence it may be worth taking a step back, including for example interaction or higher-order terms in the score estimation and testing the quality once again. If after respecification the quality indicators are still not satisfactory, it may indicate a fundamental lack of comparability of the two groups being examined. Since this is a precondition for a successful application of the matching strategy, alternative evaluation approaches should be considered (see for example the discussion in Blundell *et al.*, 2005).

It should also be noted that different matching estimators balance the covariates to different degrees. Hence, for a given estimation of the propensity score, how the different matching methods balance the covariates can be used as a criterion to choose among them (leaving efficiency considerations aside).

## *3.5 Choice-Based Sampling*

An additional problem arising in evaluation studies is that samples used are often choice-based (Smith and Todd, 2005). This is a situation where programme participants are oversampled relative to their frequency in the population of eligible

persons. This type of sampling design is frequently chosen in evaluation studies to reduce the costs of data collection and to get a larger number of treated individuals (Heckman and Todd, 2004). We discuss this point briefly and suggest one correction mechanism introduced by Heckman and Todd (2004). First of all, note that under choice-based sampling weights are required to consistently estimate the probability of programme participation. Since population weights are not known in most choice-based datasets used in evaluation analysis the propensity score cannot be consistently estimated (Heckman and Todd, 2004). However, Heckman and Todd (2004) show that even with population weights unknown, matching methods can still be applied. This is the case because the odds ratio estimated using the incorrect weights (those that ignore the fact of choice-based samples) is a scalar multiple of the true odds ratio, which is itself a monotonic transformation of propensity scores. Hence, matching can be done on the (misweighted) estimate of the odds ratio (or of the log odds ratio). Clearly, with single NN matching it does not matter whether matching is performed on the odds ratio or the estimated propensity score (with wrong weights), since ranking of the observations is identical and therefore the same neighbours will be selected. However, for methods that take account of the absolute distance between observations, e.g. KM, it does matter (Smith and Todd, 2005).

### 3.6 *When to Compare and Locking-in Effects*

An important decision which has to be made in the empirical analysis is when to measure the effects. The major goal is to ensure that participants and nonparticipants are compared in the same economic environment and the same individual lifecycle position. For example, when evaluating labour market policies one possible problem which has to be taken into account is the occurrence of locking-in effects. The literature is dominated by two approaches, comparing the individuals either from the beginning of the programme or after the end of the programme. To give an example let us assume that a programme starts in January and ends in June. The latter of the two alternatives implies that the outcome of participants who reenter the labour market in July is compared with matched nonparticipants in July. There are two shortcomings to this approach. First, if the exits of participants are spread over a longer time period, it might be the case that very different economic situations are compared. Second, a further problem which arises with this approach is that it entails an endogeneity problem (Gerfin and Lechner, 2002), since the abortion of the programme may be caused by several factors which are usually not observed by the researcher.<sup>17</sup>

The above mentioned second approach is predominant in the recent evaluation literature (see e.g. Gerfin and Lechner, 2002; Sianesi, 2004) and measures the effects from the beginning of the programme. One major argument to do so concerns the policy relevance. In the above example the policy maker is faced with the decision to put an individual in January in a programme or not. He will be interested in the effect of his decision on the outcome of the participating individual in contrast with the situation if the individual would not have participated. Therefore comparing both outcomes from the beginning of the programme is a reasonable approach.

What should be kept in mind, however, is the possible occurrence of locking-in effects for the group of participants. Since they are involved in the programme, they do not have the same time to search for a new job as nonparticipants. The net effect of a programme consists of two opposite effects. First, the increased employment probability through the programme, and second, the reduced search intensity.<sup>18</sup> Since the two effects cannot be disentangled, we only observe the net effect and have to take this into account when interpreting the results. As to the fall in the search intensity, we should expect an initial negative effect from any kind of participation in a programme. However, a successful programme should overcompensate for this initial fall. So, if we are able to observe the outcome of the individuals for a reasonable time after the beginning/end of the programme, the occurrence of locking-in effects poses fewer problems but nevertheless has to be taken into account in the interpretation.

### 3.7 *Estimating the Variance of Treatment Effects*

Testing the statistical significance of treatment effects and computing their standard errors is not a straightforward thing to do. The problem is that the estimated variance of the treatment effect should also include the variance due to the estimation of the propensity score, the imputation of the common support and in the case of matching without replacement also the order in which the treated individuals are matched. These estimation steps add variation beyond the normal sampling variation (see the discussion in Heckman *et al.*, 1998b). For example, in the case of NN matching with one NN, treating the matched observations as given will understate the standard errors (Smith, 2000). Things get more complicated, since a much discussed topic in the recent evaluation literature centres around efficiency bounds of the different approaches and how to reach them. The aim of this section is to provide a brief overview of this ongoing discussion and more importantly to describe three approaches for the estimation of standard errors which are frequently used in the empirical literature.

#### *Efficiency and Large Sample Properties of Matching Estimators*

The asymptotic properties of matching and weighting estimators have been studied by for example Hahn (1998), Heckman *et al.* (1998b) and Abadie and Imbens (2006a). The results from Hahn (1998) are a good starting point for the efficiency discussion. He derives the semi-parametric efficiency bounds for ATE and ATT under various assumptions. He especially takes into account cases where the propensity score is known and where it has to be estimated.<sup>19</sup> Under the unconfoundedness assumption the asymptotic variance bounds for ATE and ATT are given by

$$\text{Var}_{ATE} = E \left[ \frac{\sigma_1^2(X)}{P(X)} + \frac{\sigma_0^2(X)}{1 - P(X)} + (E(Y(1)|X) - E(Y(0)|X) - \tau_{ATE})^2 \right] \quad (12)$$

and

$$\begin{aligned} \text{Var}_{ATT}^{P\text{Sunknown}} = E \left[ \frac{P(X)\sigma_1^2(X)}{E[P(X)]^2} + \frac{P(X)^2\sigma_0^2(X)}{E[P(X)]^2(1-P(X))} \right. \\ \left. + \frac{(E(Y(1)|X) - E(Y(0)|X) - \tau_{ATT})^2 P(X)}{E[P(X)]^2} \right] \quad (13) \end{aligned}$$

where  $\sigma_D^2(X)$  are the conditional outcome variances for treated ( $D = 1$ ) and untreated ( $D = 0$ ) observations.

There is an ongoing discussion in the literature on how the efficiency bounds are achieved and if the propensity score should be used for estimation of ATT and ATE or not. In the above cited paper Hahn (1998) shows that when using nonparametric series regression, adjusting for all covariates can achieve the efficiency bound, whereas adjusting for the propensity score does not. Hirano *et al.* (2003) show that weighting with the inverse of a nonparametric estimate of the propensity score can achieve the efficiency bound, too. Angrist and Hahn (2004) use the results from Hahn (1998) as a starting point for their analysis and note that conventional asymptotic arguments would appear to offer no justification for anything other than full control for covariates in estimation of ATEs. However, they argue that conventional asymptotic results can be misleading and provide poor guidance for researchers who face a finite sample. They develop an alternative theory and propose a panel-style estimator which can provide finite-sample efficiency gains over covariate and propensity score matching.

Heckman *et al.* (1998b) analyse large sample properties of LPM estimators for the estimation of ATT. They show that these estimators are  $\sqrt{n}$ -consistent and asymptotically normally distributed. This holds true when matching with respect to  $X$ , the known propensity score or the estimated propensity score. They conclude that none of the approaches dominates the others *per se*. In the case of matching on the known propensity score, the asymptotic variance of  $\text{Var}_{ATT}$  is not necessarily smaller than that when matching on  $X$ .<sup>20</sup>

Abadie and Imbens (2006a) analyse the asymptotic efficiency of  $n$  nearest neighbour matching when  $n$  is fixed, i.e. when the number of neighbours does not grow with increasing sample size. They show that simple matching estimators include a conditional bias term of order  $O(N^{-1/k})$ , where  $k$  is the number of continuous covariates. The bias does not disappear if  $k$  equals 2 and will dominate the large sample variance if  $k$  is at least 3. Hence, these estimators do not reach the variance bounds in (12) and (13) and are inefficient. They also describe a bias correction that removes the conditional bias asymptotically, making estimators  $\sqrt{n}$ -consistent. Additionally, they suggest a new estimator for the variance that does not require consistent nonparametric estimation of unknown functions (we will present that approach further below). Imbens (2004) highlights some caveats of these results. First, it is important to make clear that only continuous covariates should be counted in dimension  $k$ , since with discrete covariates the matching will be exact in large samples. Second, if only treated individuals are matched and the number of potential controls is much larger than the number of treated individuals, it can be justified

to ignore the bias by appealing to an asymptotic sequence where the number of potential controls increases faster than the number of treated individuals.

### *Three Approaches for the Variance Estimation*

There are a number of ways to estimate the variance of average treatment effects as displayed in equations (12) and (13). One is by ‘brute force’ (Imbens, 2004), i.e. by estimating the five components of the variance  $\sigma_0^2(X)$ ,  $\sigma_1^2(X)$ ,  $E(Y(1)|X)$ ,  $E(Y(0)|X)$  and  $P(X)$  using kernel methods or series. Even though this is consistently possible and hence the asymptotic variance will be consistent, too, Imbens (2004) notes that this might be an additional computational burden. Hence, practical alternatives are called for and we are going to present three of them. Two of them, bootstrapping and the variance approximation by Lechner (2001a), are very common in the applied literature. Additionally, we are going to present a new method from Abadie and Imbens (2006a) that is based on the distinction between average treatment effects and sample average treatment effects.

### *Bootstrapping*

One way to deal with this problem is to use bootstrapping as suggested for example by Lechner (2002). This method is a popular way to estimate standard errors in case analytical estimates are biased or unavailable.<sup>21</sup> Even though Imbens (2004) notes that there is little formal evidence to justify bootstrapping and Abadie and Imbens (2006b) even show that the standard bootstrap fails for the case of NN matching with replacement on a continuous covariate it is widely applied. An early example of use can be found in Heckman *et al.* (1997a) who report bootstrap standard errors for LLM estimators. Other application examples for bootstrapping are for example Black and Smith (2004) for NN and KM estimators or Sianesi (2004) in the context of caliper matching. Each bootstrap draw includes the reestimation of the results, including the first steps of the estimation (propensity score, common support, etc.). Repeating the bootstrapping  $R$  times leads to  $R$  bootstrap samples and  $R$  estimated average treatment effects. The distribution of these means approximates the sampling distribution (and thus the standard error) of the population mean. Clearly, one practical problem arises because bootstrapping is very time consuming and might therefore not be feasible in some cases.

### *Variance Approximation by Lechner:*

An alternative is suggested by Lechner (2001a). For the estimated ATT obtained via NN matching the following formula applies:

$$\text{Var}(\hat{\tau}_{ATT}) = \frac{1}{N_1} \text{Var}(Y(1) | D = 1) + \frac{\sum_{j \in \{D=0\}} (w_j)^2}{(N_1)^2} \cdot \text{Var}(Y(0) | D = 0) \quad (14)$$

where  $N_1$  is the number of matched treated individuals and  $w_j$  is the number of times individual  $j$  from the control group has been used, i.e. this takes into account

that matching is performed with replacement. If no unit is matched more than once, the formula coincides with the ‘usual’ variance formula. By using this formula to estimate the variance of the treatment effect at time  $t$ , we assume independent observations and fixed weights. Furthermore we assume homoscedasticity of the variances of the outcome variables within treatment and control group and that the outcome variances do not depend on the estimated propensity score. This approach can be justified by simulation results from Lechner (2002) who finds little difference between bootstrapped variances and the variances calculated according to equation (14).

### *Variance Estimators by Abadie and Imbens*

To introduce this variance estimator, some additional notation is needed. Abadie and Imbens (2006a) explicitly distinguish average treatment effects given in Section 2 from sample average treatment effects. The latter estimators focus on the average treatment effects in the specific sample rather than in the population at large. Hence, the sample average treatment effect for the treated (SATT) is given by

$$\tau_{SATT} = \frac{1}{N_1} \sum_{i \in \{D=1\}} [Y_i(1) - Y_i(0)] \quad (15)$$

Abadie and Imbens (2006a) derived a matching variance estimator that does not require additional nonparametric estimation. The basic idea is that even though the asymptotic variance depends on the conditional variances  $\sigma_1^2(X)$  and  $\sigma_0^2(X)$ , one actually need not estimate these variances consistently at all values of the covariates. Instead only the average of this variance over the distribution weighted by the inverse of  $P(X)$  and  $1 - P(X)$  is needed. The variance of SATT can then be estimated by

$$\text{Var}_{SATT} = \frac{1}{N_1} \sum_{i=1}^N \left( D_i - (1 - D_i) \cdot \frac{K_M(i)}{M} \right)^2 \hat{\sigma}_{D_i}^2(X_i) \quad (16)$$

where  $M$  is the number of matches and  $K_M(i)$  is the number of times unit  $i$  is used as a match.

It should be noted that the estimation of the conditional variances requires estimation of conditional outcome variances  $\sigma_D^2(X_i)$ . Abadie and Imbens (2006a) offer two options. With the first option one assumes that the treatment effect is constant for all individuals  $i$  and that  $\sigma_D^2(X_i)$  does not vary with  $X$  or  $D$ . This is the assumption of homoscedasticity, whereas heteroscedasticity is allowed in the second option, where it is explicitly allowed that  $\sigma_D^2(X_i)$  differ in  $D$  and  $X$ .<sup>22</sup>

### *3.8 Combined and Other Propensity Score Methods*

What we have discussed so far is the estimation of treatment effects under unconfoundedness with (propensity score) matching estimators. Imbens (2004) notes that one evaluation method alone is often sufficient to obtain consistent or even efficient estimates. However, combining evaluation methods is a straightforward



way to improve their performance by eliminating remaining bias and/or improving precision. In this section we address three combined methods.

First, we introduce an estimator which combines matching with the DID approach. By doing so, a possible bias due to time-invariant unobservables is eliminated. Second, we present a regression-adjusted matching estimator that combines matching with regression. This can be useful because matching does not address the relation between covariates and outcome. Additionally, if covariates appear seriously imbalanced after propensity score matching (inexact or imperfect matching) a bias-correction procedure after matching may help to improve estimates. Third, we present how weighting on the propensity score can be used to obtain a balanced sample of treated and untreated individuals.<sup>23</sup>

### *Conditional DID or DID Matching Estimator*

The matching estimator described so far assumes that after conditioning on a set of observable characteristics, (mean) outcomes are independent of programme participation. The conditional DID or DID matching estimator relaxes this assumption and allows for unobservable but temporally invariant differences in outcomes between participants and nonparticipants. This is done by comparing the conditional before–after outcome of participants with those of nonparticipants. DID matching was first suggested by Heckman *et al.* (1998a). It extends the conventional DID estimator by defining outcomes conditional on the propensity score and using semiparametric methods to construct the differences. Therefore it is superior to DID as it does not impose linear functional form restrictions in estimating the conditional expectation of the outcome variable and it reweights the observations according to the weighting function of the matching estimator (Smith and Todd, 2005). If the parameter of interest is ATT, the DID propensity score matching estimator is based on the following identifying assumption:

$$E[Y_t(0) - Y_{t'}(0)|P(X), D = 1] = E[Y_t(0) - Y_{t'}(0)|P(X), D = 0] \quad (17)$$

where ( $t$ ) is the post- and ( $t'$ ) is the pretreatment period. It also requires the common support condition to hold. If panel data on participants and nonparticipants are available, it can be easily implemented by constructing propensity score matching estimates in both periods and calculating the difference between them.<sup>24</sup> Smith and Todd (2005) find that such estimators are more robust than traditional cross-section matching estimators.

### *Regression-Adjusted and Bias-Corrected Matching Estimators*

The regression-adjusted matching estimator (developed by Heckman *et al.*, 1997a, 1998a) combines LLM on the propensity score with regression adjustment on covariates. By utilizing information on the functional form of outcome equations and by incorporating exclusion restrictions across outcome and participation equation, it extends classical matching methods. Heckman *et al.* (1998b) present a proof of consistency and asymptotic normality of this estimator. Navarro-Lozano (2002)

provides a nice example for an application by evaluating a popular training programme in Mexico.

In cases where (substantial) differences in covariates between matched pairs remain after matching, additional regression adjustments may be helpful to reduce such differences. If matching is not exact, there will be some discrepancies that lead to a potential bias. The basic idea of the bias-correction estimators is to use the difference in the covariates to reduce the bias of the matching estimator. Rubin (1973, 1979) first proposed several matched sample regression adjustments in the context of Mahalanobis metric matching and they have been more recently discussed by Abadie and Imbens (2006a) and Imbens (2004).

### *Weighting on the Propensity Score*

Imbens (2004) notes that propensity scores can also be used as weights to obtain a balanced sample of treated and untreated individuals.<sup>25</sup> Such estimators can be written as the difference between a weighted average of the outcomes for the treated and untreated individuals, where units are weighted by the reciprocal of the probability of receiving treatment.<sup>26</sup> An unattractive feature of such estimators is that the weights do not necessarily add up to one. One approach to improve the propensity score weighting estimator is to normalize the weights to unity. If the propensity score is known, the estimator can directly be implemented. But, even in randomized settings where the propensity score is known, Hirano *et al.* (2003) show that it could be advantageous in terms of efficiency considerations to use the estimated rather than the 'true' propensity score. However, as Zhao (2004) notes, the way propensity scores are estimated is crucial when implementing weighting estimators and mis-specification of the propensity score may lead to substantial bias.<sup>27</sup>

### *3.9 Sensitivity Analysis*

Checking the sensitivity of the estimated results becomes an increasingly important topic in the applied evaluation literature. We will address two possible topics for a sensitivity analysis in this section. First, we are going to discuss approaches that allow the researcher to assess the sensitivity of the results with respect to deviations from the identifying assumption. Second, we show how to incorporate information from those individuals who failed the common support restriction to calculate bounds of the parameter of interest (if all individuals from the sample at hand would have been included).

### *Deviations from Unconfoundedness or Unobserved Heterogeneity*

We have outlined in Section 2 that the estimation of treatment effects with matching estimators is based on the unconfoundedness or selection on observables assumption. However, if there are unobserved variables which affect assignment into treatment and the outcome variable simultaneously, a 'hidden bias' might arise (Rosenbaum,

2002). It should be clear that matching estimators are not robust against this 'hidden bias'. Researchers become increasingly aware that it is important to test the robustness of results to departures from the identifying assumption. Since it is not possible to estimate the magnitude of selection bias with nonexperimental data, the problem can be addressed by sensitivity analysis. Even though the idea for such analyses reaches far back in the literature only a few applied studies take them into account. However, it seems that this topic has come back into the mind of applied researchers and will become more important in the next few years. The aim of this section is to give a brief overview of some of the suggested methods.<sup>28</sup>

One of the earliest examples for sensitivity analysis in the evaluation context can be found in Rosenbaum and Rubin (1983a). They propose to assess the sensitivity of ATE with respect to assumptions about an unobserved binary covariate that is associated both with the treatment and the response. The basic idea is that treatment is not unconfounded given the set of observable characteristics  $X$  but would be unconfounded given  $X$  and an unobservable covariate  $U$ . Based on different sets of assumptions about the distribution of  $U$  and its association with  $D$  and the outcomes  $Y(0)$  and  $Y(1)$  it is then possible to check the sensitivity of the results with respect to variations in these assumptions.

Imbens (2003) builds on this approach but does not formulate the sensitivity in terms of coefficients on the unobserved covariate and rather presents the sensitivity results in terms of partial  $R^2$ s. This eases the interpretation and additionally allows a comparison of the partial  $R^2$ s of the unobserved covariates to those for the observed covariates in order to facilitate judgements regarding the plausibility of values necessary to substantially change results obtained under exogeneity. Both approaches use a parametric model as the basis for estimating ATEs. Parametrization is not needed, however, in the following two approaches.

The first approach was proposed by Rosenbaum (2002) and has been recently applied in Aakvik (2001), DiPrete and Gangl (2004) and Caliendo *et al.* (2007). The basic question to be answered here is whether inference about treatment effects may be altered by unobserved factors. In other words, one wants to determine how strongly an unmeasured variable must influence the selection process in order to undermine the implications of matching analysis. To do so it is assumed that the participation probability  $\pi_i$  is not only determined by observable factors ( $x_i$ ) but also by an unobservable component ( $u_i$ ):  $\pi_i = \Pr(D_i = 1 \mid x_i) = F(\beta x_i + \gamma u_i)$ .  $\gamma$  is the effect of  $u_i$  on the participation decision. Clearly, if the study is free of hidden bias,  $\gamma$  will be zero and the participation probability will solely be determined by  $x_i$ . However, if there is hidden bias, two individuals with the same observed covariates  $x$  have differing chances of receiving treatment. Varying the value of  $\gamma$  allows the researcher to assess the sensitivity of the results with respect to 'hidden bias'. Based on that, bounds for significance levels and confidence intervals can be derived. (For details see Rosenbaum (2002) and Aakvik (2001). Becker and Caliendo (2007) provide an implementation in Stata).

A different approach was recently proposed by Ichino *et al.* (2006). It additionally allows assessment of the sensitivity of point estimates and specifically the sensitivity of ATT matching estimates. They derive point estimates of the ATT under different

possible scenarios of deviation from unconfoundedness. To do so they impose values of the parameters that characterize the distribution of  $U$ . Given these parameters, the value of the confounding factor for each treated and control subject is predicted and the ATT is reestimated now including the influence of the simulated  $U$  in the set of matching variables. By changing the assumptions about the distribution of  $U$ , they can assess the robustness of the ATT with respect to different hypotheses on the nature of the confounding factor. Their approach also allows one to verify whether there exists a set of plausible assumptions on  $U$  under which the estimated ATT would be driven to zero by the inclusion of  $U$  in the matching set. By modelling the nature of  $U$  based on already existing variables in the data, it is possible to assess the robustness of the estimates with respect to deviations from unconfoundedness that would occur if observed factors were omitted from the matching set.

A somewhat different strategy is to focus on estimating the causal effect of a treatment that is known to have a zero effect, e.g. by relying on the presence of multiple control groups (see the discussion in Imbens (2004) for details). If one has a group of eligible and ineligible nonparticipants, the 'treatment effect' which is known to be zero can be estimated using only the two control groups (where the 'treatment' indicator then has to be a dummy for belonging in one of the two groups). Any nonzero effect implies that at least one of the control groups is invalid. However, as Imbens (2004) points out, not rejecting the test does not imply that the unconfoundedness assumption is valid, but makes it more plausible that it holds. A good example of such a comparison can be found in Heckman *et al.* (1997a).

Overall, it should be noted that none of the tests can directly justify the unconfoundedness assumption. However, they provide some scope for making the estimates more credible if the results are not sensitive to different assumptions about unobservables factors. Clearly, if the results turn out to be very sensitive the researcher might have to think about the validity of his/her identifying assumption and consider alternative strategies. In any case, these tests should be applied more frequently.

### *Failure of Common Support*

In Section 3.3 we have presented possible approaches to implement the common support restriction. Those individuals that fall outside the region of common support have to be disregarded. But, deleting such observations yields an estimate that is only consistent for the subpopulation within the common support. However, information from those outside the common support could be useful and informative especially if treatment effects are heterogeneous.

Lechner (2001b) describes an approach to check the robustness of estimated treatment effects due to failure of common support. He incorporates information from those individuals who failed the common support restriction to calculate nonparametric bounds of the parameter of interest, if all individuals from the sample at hand would have been included. To introduce his approach some additional notation is needed. Define the population of interest with  $\Omega$  which is some subset from the space defined by treatment status ( $D = 1$  or  $D = 0$ ) and a set of covariates

$X$ .  $\Omega^{ATT}$  is defined by  $\{(D = 1) \times X\}$  and  $W^{ATT}$  is a binary variable which equals one if an observation belongs to  $\Omega^{ATT}$ . Identification of the effect is desired for  $\tau_{ATT}(\Omega^{ATT})$ . Due to missing common support the effect can only be estimated for  $\tau_{ATT}(\Omega^{ATT*})$ . This is the effect ignoring individuals from the treatment group without a comparable match. Observations within common support are denoted by the binary variable  $W^{ATT*}$  equal to one. The subset for whom such effect is not identified is  $\tilde{\Omega}^{ATT}$ .

Let  $\Pr(W^{ATT*} = 1 | W^{ATT} = 1)$  denote the share of participants within common support relative to the total number of participants and  $\lambda_0^1$  be the mean of  $Y(1)$  for individuals from the treatment group outside common support. Assume that the share of participants within common support relative to the total number of participants as well as ATT for those within the common support and  $\lambda_0^1$  are identified. Additionally, assume that the potential outcome  $Y(0)$  is bounded:  $\Pr(\underline{Y} \leq Y(0) \leq \bar{Y} | W^{ATT*} = 0, W^{ATT} = 1) = 1$ .<sup>29</sup> Given these assumptions, the bounds for ATT  $\tau_{ATT}(\Omega^{ATT}) \in [\underline{\tau}_{ATT}(\Omega^{ATT}), \bar{\tau}_{ATT}(\Omega^{ATT})]$  can be written as

$$\begin{aligned} \underline{\tau}_{ATT}(\Omega^{ATT}) &= \tau_{ATT}(\Omega^{ATT*})\Pr(W^{ATT*} = 1 | W^{ATT} = 1) \\ &\quad + (\lambda_0^1 - \bar{Y})[1 - \Pr(W^{ATT*} = 1 | W^{ATT} = 1)] \end{aligned} \quad (18)$$

$$\begin{aligned} \bar{\tau}_{ATT}(\Omega^{ATT}) &= \tau_{ATT}(\Omega^{ATT*})\Pr(W^{ATT*} = 1 | W^{ATT} = 1) \\ &\quad + (\lambda_0^1 - \underline{Y})[1 - \Pr(W^{ATT*} = 1 | W^{ATT} = 1)] \end{aligned} \quad (19)$$

Lechner (2001b) states that either ignoring the common support problem or estimating ATT only for the subpopulation within the common support can both be misleading. He recommends to routinely compute bounds analysis in order to assess the sensitivity of estimated treatment effects with respect to the common support problem and its impact on the inference drawn from subgroup estimates.

### 3.10 More Practical Issues and Recent Developments

Before we conclude the paper in the next section, we will point out some additional topics which might be of relevance in applied research. What we have discussed so far is basically a static and binary evaluation framework where an individual can participate in one programme (or not). However, in most realistic evaluation settings this framework might not be appropriate, e.g. when evaluating the effects of labour market policies. First of all, researchers are usually not confronted with only one, but a different set of programmes (programme heterogeneity). Second, an unemployed can successively enter into different programmes as long as (s)he is unemployed. Finally, choosing the right control group and the problem of random programme starts is a recently much discussed topic in the evaluation literature, too. These issues as well as a short listing of available software tools to implement matching are discussed in this section.

### *Programme Heterogeneity*

The standard evaluation framework as presented in Section 2 considers only two possible states for each individual, i.e. participation and nonparticipation. To account for programme heterogeneity, this approach has been extended by Imbens (2000) and Lechner (2001a) to the multiple treatment framework which considers the case of  $L + 1$  mutually different and exclusive treatments. For every individual only one component of the  $L + 1$  different outcomes  $\{Y(0), Y(1), \dots, Y(L)\}$  can be observed, leaving  $L$  as counterfactuals. Participation in treatment  $l$  is indicated by  $D \in \{0, 1, \dots, L\}$ . The interest lies in the causal effect of one treatment relative to another treatment on an outcome variable. Even though Lechner (2001a) defines several parameters of interest, we will focus once again on the ATT. In the multiple treatment notation, that effect is defined as a pairwise comparison of the effects of the treatments  $m$  and  $l$  for an individual randomly drawn from the group of participants in  $m$  only:

$$\tau_{ATT}^{ml} = E[Y(m) - Y(l) \mid D = m] = E[Y(m) \mid D = m] - E[Y(l) \mid D = m] \quad (20)$$

As discussed in Section 2, the causal treatment effect in the presented framework is not identified. To overcome the counterfactual situation, the unconfoundedness assumption has to be adapted to the multiple treatment framework:

$$Y(0), Y(1), \dots, Y(L) \perp\!\!\!\perp D \mid X \quad (21)$$

This assumption can be weakened when one is interested in pairwise programme comparisons only. If we further assume that those receiving treatment  $m$  have a counterpart in the comparison group, i.e. if there is common support, the counterfactual mean can be constructed as  $E[Y(l) \mid D = m, X]$ . Lechner (2001a) also shows that the generalization of the balancing property holds for the case of multiple treatments as well. To estimate  $\tau_{ATT}^{ml}$  matching can be done by using the conditional choice probability of treatment  $m$  given either treatment  $m$  or  $l$  and covariates  $X$  as a balancing score:

$$P(D = m \mid X, D \in \{m, l\}) = \frac{P(D = m \mid X)}{P(D = m \mid X) + P(D = l \mid X)} \quad (22)$$

If the conditional choice probability is modelled directly, no information from subsamples other than those containing participants in  $m$  and  $l$  is needed and one is basically back in the binary treatment framework. Since the choice probabilities will not be known *a priori*, they have to be replaced by an estimate, e.g. a probit model. If all values of  $m$  and  $l$  are of interest, the whole sample is needed for identification. In that case either the binary conditional probabilities can be estimated or a structural approach can be used where a complete choice problem is formulated in one model and estimated on the full sample, e.g. with a multinomial probit model. We have discussed the (dis-)advantages of the multinomial modelling in comparison to discrete estimation of binomial models already in Section 3.1.

### *Sequential Matching Estimators*

Extending the standard evaluation framework for the case where individuals can participate in subsequent treatments has been recently proposed by Lechner and Miquel (2005).<sup>30</sup> These ‘programme careers’ cannot be addressed properly in the basic framework. Problems occur because the assignment into a subsequent programme is not independent of the assignment into previous programmes. Additionally, outcomes in subsequent periods will be influenced by previous participation decisions. Hence, a dynamic selection problem arises. Most empirical work about dynamic selection problems ignores intermediate outcomes and treats the sequence participation as being determined from the start. Mainly, problems are circumvented by either estimating the effect of the first programme only (see e.g. Gerfin and Lechner, 2002) or applying the static framework subsequently (see e.g. Bergemann *et al.*, 2001). The sequential matching framework is a powerful tool and is applicable for situations where individuals can participate more than once in a programme and where it is possible to identify treatment sequences. It allows intermediate outcomes to play a role in the participation decision for sequential participation and thus allows estimation in a dynamic context. To our knowledge Lechner (2004) is the only application so far and hence practical experiences with sequential matching estimators are rather limited.

### *Choosing the Right Control Group – Random Programme Starts*

Another important topic in applied evaluation research is to choose an appropriate control group. In the ‘usual’ evaluation set-up for matching estimators, we have a group of participants and a group of nonparticipants. Both groups are usually observed from a certain starting point  $t$  to an end point  $T$ . The researcher does not have any information outside this limited time interval. Controls are defined as those individuals who did not participate in any programme in  $[t, T]$ , whereas participants are those individuals who took part in a programme for a certain interval  $\tau$  in  $[t, T]$ .

In a series of papers, Sianesi (2001, 2004) casts doubt if this standard approach is appropriate. She suggests a solution which is based on a redefinition of the control group. Instead of defining controls as those who never participate, she defines controls as those who did not participate until a certain time period. Hence, the corresponding parameter of interest in this setting is then defined as the effect of joining a programme now in contrast to waiting longer. Fredriksson and Johansson (2004) formalize her approach and argue that the standard way of defining a control group might lead to biased results, because the unconfoundedness assumption might be violated. The reason for this is that in the standard approach the treatment indicator itself is defined conditional on future outcomes. In fact, in the context of labour market policies it can be argued that an unemployed individual will join a programme at some time, provided his unemployment spell is long enough (Sianesi, 2004). Hence, if the reason for nonparticipation is that the individual has found a job before a participation in the programme was offered or considered, it leads to negatively biased effects.

### *Available Software to Implement Matching*

The bulk of software tools to implement matching and estimate treatment effects is growing and allows researchers to choose the appropriate tool for their purposes. The most commonly used platform for these tools is Stata and we will present the three most distributed ones here. Becker and Ichino (2002) provide a programme for PSM estimators (*pscore*, *attnd*, *attnw*, *attr*, *atts*, *attk*) which includes estimation routines for NN, kernel, radius, and stratification matching. To obtain standard errors the user can choose between bootstrapping and the variance approximation proposed by Lechner (2001a). Additionally the authors offer balancing tests (blocking, stratification) as discussed in Section 3.4.

Leuven and Sianesi (2003) provide the programme *psmatch2* for implementing different kinds of matching estimators including covariate and propensity score matching. It includes NN and caliper matching (with and without replacement), KM, radius matching, LLM and Mahalanobis metric (covariate) matching. Furthermore, this programme includes routines for common support graphing (*psgraph*) and covariate imbalance testing (*pstest*). Standard errors are obtained using bootstrapping methods.

Finally, Abadie *et al.* (2004) offer the programme *nnmatch* for implementing covariate matching, where the user can choose between several different distance metrics. Variance approximations as proposed by Abadie and Imbens (2006a) are implemented to obtain standard errors of treatment effects.

## **4. Conclusion**

The aim of this paper was to give some guidance for the implementation of propensity score matching. Basically five implementation steps have to be considered when using PSM (as depicted in Figure 1). The discussion has made clear that a researcher faces a lot of decisions during implementation and that it is not always an easy task to give recommendations for a certain approach. Table 2 summarizes the main findings of this paper and also highlights sections where information for each implementation step can be found.

The first step of implementation is the estimation of the propensity score. We have shown that the choice of the underlying model is relatively unproblematic in the binary case whereas for the multiple treatment case one should either use a multinomial probit model or a series of binary probits (logits). After having decided about which model to be used, the next question concerns the variables to be included in the model. We have argued that the decision should be based on economic theory, a sound knowledge of previous research and also information about the institutional settings. We have also presented several statistical strategies which may help to determine the choice. If it is felt that some variables play a specifically important role in determining participation and outcomes, one can use an 'overweighting' strategy, for example by carrying out matching on subpopulations.

The second implementation step is the choice among different matching algorithms. We have argued that there is no algorithm which dominates in all data



**Table 2.** Implementation of Propensity Score Matching.

Step	Decisions, questions and solutions	Section
<b>1. Estimation of propensity score</b>		
Model choice	<ul style="list-style-type: none"> <li>◊ Unproblematic in the binary treatment case (logit/probit)</li> <li>◊ In the multiple treatment case multinomial probit or series of binomial models should be preferred</li> <li>◊ Variables should not be influenced by participation (or anticipation) and must satisfy CIA</li> </ul>	3.1 3.1
Variable choice	<ul style="list-style-type: none"> <li>◊ Choose variables by economic theory and previous empirical evidence</li> <li>◊ ‘Hit or miss’ method, stepwise augmentation, leave-one-out cross-validation</li> <li>◊ ‘Overweighting’ by matching on subpopulations or insisting on perfect match</li> </ul>	3.1 3.1 3.1
<b>2. Choice among alternative matching algorithms</b>		
Matching algorithms	<ul style="list-style-type: none"> <li>◊ The choice (e.g. NN matching with or without replacement, caliper or kernel matching) depends on the sample size, the available number of treated/control observations and the distribution of the estimated propensity score</li> <li>→ Trade-offs between bias and efficiency!</li> </ul>	3.2
<b>3. Check overlap and common support</b>		
Common support	<ul style="list-style-type: none"> <li>◊ Treatment effects can be estimated only over the CS region!</li> </ul>	3.3
→ Tests	Visual analysis of propensity score distributions	3.3
→ Implementation	‘Minima and maxima comparison’ or ‘trimming’ method Alternative: Caliper matching	3.3
<b>4.1 Assessing the matching quality</b>		
Balancing property	<ul style="list-style-type: none"> <li>◊ Is the matching procedure able to balance the distribution of relevant covariates?</li> <li>◊ If matching was not successful go back to step 1 and include higher-order terms, interaction variables or different covariates</li> </ul>	3.4 ↔ Step 1

Table 2. *Continued*

Step	Decisions, questions and solutions	Section
→ Tests	<ul style="list-style-type: none"> <li>◇ After that, if matching is still not successful it may indicate a fundamental lack of comparability between treatment and control group</li> <li>→ Consider alternative evaluation approaches</li> <li>Standardized bias, <math>t</math>-test, stratification test, joint significance and pseudo-<math>R^2</math></li> </ul>	3.4
<b>4.2 Calculation of treatment effects</b>		
Choice-based sample	◇ Sample is choice-based? Match on the odds ratio instead of the propensity score	3.5
When to compare	◇ Compare from the beginning of the programme to avoid endogeneity problems	3.6
	→ Pay attention to the possible occurrence of locking-in effects	3.6
Standard errors	◇ Calculate standard errors by bootstrapping or variance approximation	3.7
Combined methods	◇ Think about combining PSM with other evaluation methods to possibly eliminate remaining bias and/or improve precision	3.8
<b>5. Sensitivity analysis</b>		
Hidden bias	◇ Test the sensitivity of estimated treatment effects with respect to unobserved covariates	3.9
	→ If results are very sensitive reconsider identifying assumption and consider alternative estimators	
Common support	◇ Test the sensitivity of estimated treatment effects with respect to the common support problem	3.9
	→ Calculate Lechner bounds. If results are very sensitive reconsider variable choice	↔ Step 1

CS, common support; NN, nearest neighbour; CIA, conditional independence assumption.

situations and that the choice involves a trade-off between bias and efficiency. The performance of different matching algorithms varies case-by-case and depends largely on the data sample. If results among different algorithms differ substantially, further investigations may be needed to reveal the source of disparity.

The discussion has also emphasized that treatment effects can only be estimated in the region of common support. To identify this region we recommend to start with a visual analysis of the propensity score distributions in the treatment and comparison group. Based on that, different strategies can be applied to implement the common support condition, e.g. by 'minima and maxima comparison' or 'trimming', where the latter approach has some advantages when observations are close to the 'minima and maxima' bounds and if the density in the tails of the distribution is very thin.

Since we do not condition on all covariates but on the propensity score we have to check in the next step if the matching procedure is able to balance the distribution of these covariates in the treatment and comparison group. We have presented several procedures to do so, including SB, *t*-test, stratification test, joint significance and pseudo- $R^2$ . If the quality indicators are not satisfactory, one should go back to step 1 of the implementation procedure and include higher-order or interaction terms of the existing covariates or choose different covariates (if available). If, after that, the matching quality is still not acceptable, this may indicate a lack of comparability of the two groups being examined. Since this is a precondition for a successful application of the matching estimator, one has to consider alternative evaluation approaches.

However, if the matching quality is satisfactory one can move on to estimate the treatment effects. The estimation of standard errors is a much discussed topic in the recent evaluation literature. We have briefly discussed (some) efficiency and large sample properties of matching estimators and highlighted that the discussion in this direction is not final yet. Keeping that in mind, we have introduced three approaches for the estimation of variances of treatment effects which are used, i.e. bootstrapping methods, the variance approximation proposed in Lechner (2001a) and the variance estimators proposed by Abadie and Imbens (2006a). Another important decision is 'when to measure the effects?' where we argue that it is preferable to measure the effects from the beginning of the treatment. Clearly, what has to be kept in mind for the interpretation is the possible occurrence of locking-in effects.

Finally, a last step of matching analysis is to test the sensitivity of results with respect to deviations from the identifying assumption, e.g. when there are unobserved variables which affect assignment into treatment and the outcome variable leading to a 'hidden bias'. We have pointed out that matching estimators are not robust against this bias and that researchers become increasingly aware that it is important to test the sensitivity of their results. If the results are sensitive and if the researcher has doubts about the validity of the unconfoundedness assumption he should either consider using alternative identifying assumptions or combine PSM with other evaluation approaches.

We have introduced some possible combinations in Section 3.8 where we presented the DID matching estimator, which eliminates a possible bias due to time-invariant unobservables, as well as regression-adjusted and bias-corrected matching

estimators. All approaches aim to improve the performance of the estimates by eliminating remaining bias and/or improving precision. Last, in Section 3.10 we discussed some additional topics which might be of relevance in applied research, e.g. programme heterogeneity, sequential matching estimators and the choice of the right control group.

To conclude, we have discussed several issues surrounding the implementation of PSM. We hope to give some guidance for researchers who believe that their data are strong enough to credibly justify the unconfoundedness assumption and who want to use PSM.

## Acknowledgements

The authors thank Sascha O. Becker, Arne Uhlendorff and three anonymous referees for valuable comments which helped to improve the paper. All remaining errors are our own.

## Notes

1. See e.g. Rubin (1974), Rosenbaum and Rubin (1983, 1985a) or Lechner (1998).
2. The decision whether to apply PSM or covariate matching (CVM) as well as to include the propensity score as an additional covariate into Mahalanobis metric matching will not be discussed in this paper. With CVM distance measures like the Mahalanobis distance are used to calculate similarity of two individuals in terms of covariate values and the matching is done on these distances. The interested reader is referred to Imbens (2004) or Abadie and Imbens (2006a) who develop covariate and bias-adjusted matching estimators and Zhao (2004) who discusses the basic differences between PSM and CVM.
3. Note that the stable unit treatment value assumption (SUTVA) has to be made (see Rubin (1980) or Holland (1986) for a further discussion of this concept). It requires in particular that an individual's potential outcomes depend on his own participation only and not on the treatment status of other individuals in the population. Peer-effects and general equilibrium effects are ruled out by this assumption (Sianesi, 2004).
4. For distributions of programme impacts, the interested reader is referred to Heckman *et al.* (1997b). Another parameter one might think of is the average treatment effect on the untreated (ATU):  $\tau_{ATU} = E(\tau \mid D = 0) = E[Y(1) \mid D = 0] - E[Y(0) \mid D = 0]$ . The treatment effect for those individuals who actually did not participate in the programme is typically an interesting measure for decisions about extending some treatment to a group that was formerly excluded from treatment.
5. See Smith (2000) for a discussion about advantages and disadvantages of social experiments.
6. See Heckman and Robb (1985), Heckman *et al.* (1999), Blundell and Costa Dias (2002) or Caliendo and Hujer (2006) for a broader overview of evaluation strategies including situations where selection is also based on unobservable characteristics.
7. Once again, to identify ATT it is sufficient to assume  $Y(0) \perp D \mid P(X)$ .
8. Especially the 'independence from irrelevant alternatives' assumption (IIA) is critical. It basically states that the odds ratio between two alternatives is independent of other alternatives. This assumption is convenient for estimation but not appealing from an economic or behavioural point of view (for details see e.g. Greene, 2003).

9. See e.g. Breiman *et al.* (1984) for a theoretical discussion and Heckman *et al.* (1998a) or Smith and Todd (2005) for applications.
10. See Smith and Todd (2005) or Imbens (2004) for more technical details.
11. This shortcoming is circumvented by an optimal full matching estimator which works backwards and rearranges already matched treated individuals if some specific treated individual turns out to be a better (closer) match for an untreated previously matched individual (see Gu and Rosenbaum (1993) or Augurzyk and Kluve (2007) for detailed descriptions).
12. It should be noted that the increase in the variance is due to the imposition of the common support and hence variance comparisons between matching estimators with and without caliper are not obvious.
13. The trimming method was first suggested by Heckman *et al.* (1997a, 1998a).
14. For details on how to estimate the cut-off trimming level see Smith and Todd (2005). Galdo (2004) notes that the determination of the smoothing parameter is critical here. If the distribution is skewed to the right for participants and skewed to the left for nonparticipants, assuming a normal distribution may be very misleading.
15. In a most recent paper Crump *et al.* (2005) point out that both methods presented here are somewhat informal in the sense that they rely on arbitrary choices regarding thresholds for discarding observations. They develop formal methods for addressing lack of support and especially provide new estimators based on a redefinition of the estimand.
16. Smith and Todd (2005) note that this theorem holds for any  $X$ , including those that do not satisfy the CIA required to justify matching. As such, the theorem is not informative about which set of variables to include in  $X$ .
17. It may be the case for example that a participant receives a job offer and refuses to participate because he thinks the programme is not enhancing his employment prospects or because lack of motivation. As long as the reasons for abortion are not identified, an endogeneity problem arises.
18. These ideas date back to Becker (1964) who makes the point that human capital investments are composed of an investment period, in which one incurs the opportunity cost of not working, and a payoff period, in which one's employment and/or wage are higher than they would have been without the investment.
19. Hahn (1998) shows that the propensity score does not play a role for the estimation of ATE, but knowledge of the propensity score matters for the estimation of ATT.
20. Whereas matching on  $X$  involves  $k$ -dimensional nonparametric regression function estimation (where  $k = 1, \dots, K$  are the number of covariates), matching on  $P(X)$  only involves one-dimensional nonparametric regression function estimation. Thus from the perspective of bias, matching on  $P(X)$  is preferable, since it allows  $\sqrt{n}$ -consistent estimation of  $\tau_{ATT}$  for a wider class of models (Heckman *et al.*, 1998b).
21. See Brownstone and Valletta (2001) for a discussion of bootstrapping methods.
22. See Abadie and Imbens (2006a) and Abadie *et al.* (2004) for details about the derivation of the relevant formulas and some easy implementable examples.
23. Due to space constraints we cannot address all possible combinations. For a combination of propensity score methods with an instrumental variable approach the interested reader is referred to Abadie (2003), and how to combine DID with weighting on the propensity score has been recently proposed by Abadie (2005).
24. Smith and Todd (2005) present a variant of this estimator when repeated cross-section data are used instead of panel data. With repeated cross-section data the

- identity of future participants and nonparticipants may not be known in  $t'$ , Blundell and Costa Dias (2000) suggest a solution for that case.
25. See e.g. Imbens (2004) or Wooldridge (2004), Section 18.3.2, for a formal description of weighting on propensity score estimators.
  26. See Imbens (2004) for a formal proof that this weighting estimator removes the bias due to different distributions of the covariates between treated and untreated individuals.
  27. In the recent methodological literature several estimators have been proposed that combine weighting on propensity score estimators with other methods. Due to space limitations we cannot address these topics. The interested reader is referred to for example Hirano and Imbens (2002) who apply a combined weighting on propensity score and regression adjustment estimator in their analysis or Abadie (2005) who combines DID and weighting estimators.
  28. See Ichino *et al.* (2006) or Imbens (2004) for a more detailed discussion of these topics.
  29. For example, if the outcome variable of interest is a dummy variable,  $Y(0)$  is bounded in  $[0, 1]$ .
  30. See Lechner and Miquel (2005) and Lechner (2004) for a sequential (three-periods, two-treatments) matching framework.

## References

- Aakvik, A. (2001) Bounding a matching estimator: the case of a Norwegian training program. *Oxford Bulletin of Economics and Statistics* 63(1): 115–143.
- Abadie, A. (2003) Semiparametric instrumental variable estimation of treatment response models. *Journal of Econometrics* 113: 231–263.
- Abadie, A. (2005) Semiparametric difference-in-differences estimators. *Review of Economic Studies* 72(1): 1–19.
- Abadie, A. and Imbens, G. (2006a) Large sample properties of matching estimators for average treatment effects. *Econometrica* 74(1): 235–267.
- Abadie, A. and Imbens, G. (2006b) On the failure of the bootstrap for matching estimators. Working Paper, Harvard University.
- Abadie, A., Drukker, D., Leber Herr, J. and Imbens, G. (2004) Implementing matching estimators for average treatment effects in STATA. *Stata Journal* 4(3): 290–311.
- Angrist, J. and Hahn, J. (2004) When to control for covariates? Panel-asymptotic results for estimates of treatment effects. *Review of Economics and Statistics* 86(1): 58–72.
- Augurzky, B. and Kluve, J. (2007) Assessing the performance of matching algorithms when selection into treatment is strong. *Journal of Applied Econometrics* 22(3): 533–557.
- Augurzky, B. and Schmidt, C. (2001) The propensity score: a means to an end. Discussion Paper No. 271, IZA.
- Becker, S.O. (1964) *Human Capital*. New York: Columbia University Press.
- Becker, S.O. and Ichino, A. (2002) Estimation of average treatment effects based on propensity scores. *Stata Journal* 2(4): 358–377.
- Becker, S.O. and Caliendo, M. (2007) Sensitivity analysis for average treatment effect. *Stata Journal* 7(1): 71–83.
- Bergemann, A., Fitzenberger, B. and Speckesser, S. (2001) Evaluating the employment effects of public sector sponsored training in East Germany: conditional difference-in-differences and Ashenfelters' dip. Discussion Paper, University of Mannheim.

- Black, D. and Smith, J. (2004) How robust is the evidence on the effects of the college quality? Evidence from matching. *Journal of Econometrics* 121(1): 99–124.
- Blundell, R. and Costa Dias, M. (2000) Evaluation methods for non-experimental data. *Fiscal Studies* 21(4): 427–468.
- Blundell, R. and Costa Dias, M. (2002) Alternative approaches to evaluation in empirical microeconomics. *Portuguese Economic Journal* 1: 91–115.
- Blundell, R., Dearden, L. and Sianesi, B. (2005) Evaluating the impact of education on earnings in the UK: models, methods and results from the NCDS. *Journal of the Royal Statistical Society, Series A* 168(3): 473–512.
- Brand, J.E. and Halaby, C.N. (2006) Regression and matching estimates of the effects of elite college attendance on educational and career achievement. *Social Science Research*, 35(3): 749–770.
- Breiman, L., Friedman, J., Olsen, R. and Stone, C. (1984) *Classification and Regression Trees*. Belmont, CA: Wadsworth International Group.
- Brodaty, T., Crepon, B. and Fougere, D. (2001) Using matching estimators to evaluate alternative youth employment programs: evidence from France, 1986–1988. In M. Lechner and F. Pfeiffer (eds), *Econometric Evaluation of Labour Market Policies* (pp. 85–123). Heidelberg: Physica.
- Brownstone, D. and Valletta, R. (2001) The bootstrap and multiple imputations: harnessing increased computing power for improved statistical tests. *Journal of Economic Perspectives* 15(4): 129–141.
- Bryson, A. (2002) The union membership wage premium: an analysis using propensity score matching. Discussion Paper No. 530, Centre for Economic Performance, London.
- Bryson, A., Dorsett, R. and Purdon, S. (2002) The use of propensity score matching in the evaluation of labour market policies. Working Paper No. 4, Department for Work and Pensions.
- Caliendo, M. and Hujer, R. (2006) The microeconomic estimation of treatment effects – an overview. *Allgemeines Statistisches Archiv* 90(1): 197–212.
- Caliendo, M., Hujer, R. and Thomsen, S. (2007) The employment effects of job creation schemes in Germany – a microeconomic evaluation. IZA Discussion Paper No. 1512. *Advances in Econometrics* 21, forthcoming.
- Cochran, W. (1968) The effectiveness of adjustment by subclassification in removing bias in observational studies. *Biometrics* 24: 295–314.
- Crump, R., Hotz, V., Imbens, G. and Mitnik, O. (2005) Moving the goalposts: addressing limited overlap in estimation of average treatment effects by changing the estimand. Working Paper, University of California at Berkeley.
- Davies, R. and Kim, S. (2003) Matching and the estimated impact of interlisting. Discussion Paper in Finance No. 2001-11, ISMA Centre, Reading.
- Dehejia, R. (2005) Practical propensity score matching: a reply to Smith and Todd. *Journal of Econometrics* 125: 355–364.
- Dehejia, R.H. and Wahba, S. (1999) Causal effects in nonexperimental studies: reevaluating the evaluation of training programs. *Journal of the American Statistical Association* 94(448): 1053–1062.
- Dehejia, R.H. and Wahba, S. (2002) Propensity score matching methods for nonexperimental causal studies. *Review of Economics and Statistics* 84(1): 151–161.
- DiNardo, J. and Tobias, J. (2001) Nonparametric density and regression estimation. *Journal of Economic Perspectives* 15(4): 11–28.
- DiPrete, T. and Gangl, M. (2004) Assessing bias in the estimation of causal effects: Rosenbaum bounds on matching estimators and instrumental variables estimation with imperfect instruments. *Sociological Methodology* 34: 271–310.
- Fredriksson, P. and Johansson, P. (2004) Dynamic treatment assignment – the consequences for evaluations using observational data. Discussion Paper No. 1062, IZA.

- Galdo, J. (2004) Evaluating the performance of non-experimental estimators: evidence from a randomized UI program. Working Paper, Centre for Policy Research, Toronto.
- Gerfin, M. and Lechner, M. (2002) A microeconomic evaluation of the active labour market policy in Switzerland. *The Economic Journal* 112(482): 854–893.
- Greene, W.H. (2003) *Econometric Analysis*. New York: New York University.
- Gu, X.S. and Rosenbaum, P.R. (1993) Comparison of multivariate matching methods: structures, distances and algorithms. *Journal of Computational and Graphical Statistics* 2: 405–420.
- Hahn, J. (1998) On the role of the propensity score in efficient semiparametric estimation of average treatment effects. *Econometrica* 66(2): 315–331.
- Ham, J., Li, X. and Reagan, P. (2004) Propensity score matching, a distance-based measure of migration, and the wage growth of young men. Working Paper, Department of Economics, Ohio State University.
- Heckman, J. (1997) Instrumental variables – a study of the implicit behavioral assumptions used in making program evaluations. *Journal of Human Resources* 32(3): 441–462.
- Heckman, J. and Robb, R. (1985) Alternative models for evaluating the impact of interventions. In J. Heckman and B. Singer (eds), *Longitudinal Analysis of Labor Market Data* (pp. 156–245). Cambridge: Cambridge University Press.
- Heckman, J. and Smith, J. (1999) The pre-program earnings dip and the determinants of participation in a social program: implications for simple program evaluation strategies. *Economic Journal* 109(457): 313–348.
- Heckman, J. and Todd, P. (2004) A note on adapting propensity score matching and selection models to choice based samples. Working Paper, first draft 1995, this draft Nov. 2004, University of Chicago.
- Heckman, J., Ichimura, H. and Todd, P. (1997a) Matching as an econometric evaluation estimator: evidence from evaluating a job training programme. *Review of Economic Studies* 64(4): 605–654.
- Heckman, J., Smith, J. and Clements, N. (1997b) Making the most out of programme evaluations and social experiments: accounting for heterogeneity in programme impacts. *Review of Economic Studies* 64(4): 487–535.
- Heckman, J., Ichimura, H., Smith, J. and Todd, P. (1998a) Characterizing selection bias using experimental data. *Econometrica* 66(5): 1017–1098.
- Heckman, J., Ichimura, H. and Todd, P. (1998b) Matching as an econometric evaluation estimator. *Review of Economic Studies* 65(2): 261–294.
- Heckman, J., LaLonde, R. and Smith, J. (1999) The economics and econometrics of active labor market programs. In O. Ashenfelter and D. Card (eds), *Handbook of Labor Economics*, (Vol. III, pp. 1865–2097). Amsterdam: Elsevier.
- Hirano, K. and Imbens, G. (2002) Estimation of causal effects using propensity score weighting: an application to data on right heart catheterization. *Health Services and Outcomes Research Methodology* 2(3–4): 259–278.
- Hirano, K., Imbens, G. and Ridder, G. (2003) Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica* 71(4): 1161–1189.
- Hitt, L. and Frei, F. (2002) Do better customers utilize electronic distribution channels? The case of PC banking. *Management Science* 48(6): 732–748.
- Holland, P. (1986) Statistics and causal inference. *Journal of the American Statistical Association* 81(396): 945–960.
- Ichino, A., Mealli, F. and Nannicini, T. (2006) From temporary help jobs to permanent employment: what can we learn from matching estimators and their sensitivity. Discussion Paper No. 2149, IZA, Bonn.
- Imbens, G. (2000) The role of the propensity score in estimating dose–response functions. *Biometrika* 87(3): 706–710.
- Imbens, G. (2003) Sensitivity to exogeneity assumptions in program evaluation. *American Economic Review* 93(2): 126–132.



- Imbens, G. (2004) Nonparametric estimation of average treatment effects under exogeneity: a review. *Review of Economics and Statistics* 86(1): 4–29.
- Lechner, M. (1998) Mikröökonomische Evaluationsstudien: Anmerkungen zu Theorie und Praxis. In F. Pfeiffer and W. Pohlmeier (eds), *Qualifikation, Weiterbildung und Arbeitsmarkterfolg. ZEW-Wirtschaftsanalysen Band 31*. Baden-Baden: Nomos-Verlag.
- Lechner, M. (1999) Earnings and employment effects of continuous off-the-job training in East Germany after unification. *Journal of Business Economic Statistics* 17(1): 74–90.
- Lechner, M. (2000) An evaluation of public sector sponsored continuous vocational training programs in East Germany. *Journal of Human Resources* 35(2): 347–375.
- Lechner, M. (2001a) Identification and estimation of causal effects of multiple treatments under the conditional independence assumption. In M. Lechner and F. Pfeiffer (eds), *Econometric Evaluation of Labour Market Policies* (pp. 1–18). Heidelberg: Physica.
- Lechner, M. (2001b) A note on the common support problem in applied evaluation studies. Discussion Paper No. 2001-01, University of St Gallen, SIAW.
- Lechner, M. (2002) Some practical issues in the evaluation of heterogeneous labour market programmes by matching methods. *Journal of the Royal Statistical Society, A* 165: 59–82.
- Lechner, M. (2004) Sequential matching estimation of dynamic causal models. Discussion Paper No. 1042, IZA.
- Lechner, M. and Miquel, R. (2005) Identification of the effects of dynamic treatments by sequential conditional independence assumptions. Working Paper, SIAW.
- Leuven, E. and Sianesi, B. (2003) PSMATCH2: Stata module to perform full Mahalanobis and propensity score matching, common support graphing, and covariate imbalance testing. Software, <http://ideas.repec.org/c/boc/bocode/s432001.html>.
- Navarro-Lozano, S. (2002) Matching, selection and the propensity score: evidence from training in Mexico. Working Paper, University of Chicago.
- Pagan, A. and Ullah, A. (1999) *Nonparametric Econometrics*. Cambridge: Cambridge University Press.
- Perkins, S.M., Tu, W., Underhill, M.G., Zhou, X. and Murray, M.D. (2000) The use of propensity scores in pharmacoepidemiologic research. *Pharmacoepidemiology and Drug Safety* 9(2): 93–101.
- Rosenbaum, P.R. (2002) *Observational Studies*. New York: Springer.
- Rosenbaum, P. and Rubin, D. (1983a) Assessing sensitivity to an unobserved binary covariate in an observational study with binary outcome. *Journal of the Royal Statistical Society, Series B* 45: 212–218.
- Rosenbaum, P. and Rubin, D. (1983b) The central role of the propensity score in observational studies for causal effects. *Biometrika* 70(1): 41–50.
- Rosenbaum, P. and Rubin, D. (1984) Reducing bias in observational studies using subclassification on the propensity score. *Journal of the American Statistical Association* 79: 516–524.
- Rosenbaum, P. and Rubin, D. (1985) Constructing a control group using multivariate matched sampling methods that incorporate the propensity score. *The American Statistician* 39(1): 33–38.
- Roy, A. (1951) Some thoughts on the distribution of earnings. *Oxford Economic Papers* 3(2): 135–145.
- Rubin, D. (1973) The use of matched sampling and regression adjustment to remove bias in observational studies. *Biometrics* 29(1): 185–203.
- Rubin, D. (1974) Estimating causal effects to treatments in randomised and nonrandomised studies. *Journal of Educational Psychology* 66: 688–701.
- Rubin, D. (1979) Using multivariate matched sampling and regression adjustment to control bias in observational studies. *Journal of the American Statistical Association* 74(366): 318–328.

- Rubin, D. (1980) Comment on Basu, D. – Randomization analysis of experimental data: the Fisher randomization test. *Journal of the American Statistical Association* 75(371): 591–593.
- Rubin, D.B. and Thomas, N. (1996) Matching using estimated propensity scores: relating theory to practice. *Biometrics* 52(1): 249–264.
- Sianesi, B. (2001) An evaluation of the active labour market programmes in Sweden. Working Paper No. 2001:5, IFAU – Office of Labour Market Policy Evaluation.
- Sianesi, B. (2004) An evaluation of the Swedish system of active labour market programmes in the 1990s. *Review of Economics and Statistics* 86(1): 133–155.
- Silverman, B. (1986) *Density Estimation for Statistics and Data Analysis*. London: Chapman and Hall.
- Smith, H. (1997) Matching with multiple controls to estimate treatment effects in observational studies. *Sociological Methodology* 27: 325–353.
- Smith, J. (2000) A critical survey of empirical methods for evaluating active labor market policies. *Schweizerische Zeitschrift fuer Volkswirtschaft und Statistik* 136(3): 1–22.
- Smith, J. and Todd, P. (2005) Does matching overcome LaLonde’s critique of nonexperimental estimators? *Journal of Econometrics* 125(1–2): 305–353.
- Wooldridge, J.M. (2004) *Econometric Analysis of Cross Section and Panel Data*. Boston, MA: Massachusetts Institute of Technology.
- Zhao, Z. (2000) Data issues of using matching methods to estimate treatment effects: an illustration with NSW data set. Working Paper, China Centre for Economic Research.
- Zhao, Z. (2004) Using matching to estimate treatment effects: data requirements, matching metrics, and Monte Carlo evidence. *Review of Economics and Statistics* 86(1): 91–107.