

CHAPTER 1

Introduction

1.1 ADVANTAGES OF PANEL DATA

A longitudinal, or panel, data set is one that follows a given sample of individuals over time, and thus provides multiple observations on each individual in the sample. Panel data have become widely available in both the developed and developing countries. For instance, in the U.S., two of the most prominent panel data sets are the National Longitudinal Surveys of Labor Market Experience (NLS) and the University of Michigan's Panel Study of Income Dynamics (PSID).

The NLS began in the mid-1960s. It contains five separate longitudinal data bases covering distinct segments of the labor force: men whose ages were 45 to 59 in 1966, young men 14 to 24 in 1966, women 30 to 44 in 1967, young women 14 to 24 in 1968, and youth of both sexes 14 to 21 in 1979. In 1986, the NLS expanded to include surveys of the children born to women who participated in the National Longitudinal Survey of Youth 1979. The list of variables surveyed is running into the thousands, with the emphasis on the supply side of the labor market. Table 1.1 summarizes the NLS survey groups, the sizes of the original samples, the span of years each group has been interviewed, and the current interview status of each group (for detail, see *NLS Handbook 2000*, U.S. Department of Labor, Bureau of Labor Statistics).

The PSID began with collection of annual economic information from a representative national sample of about 6,000 families and 15,000 individuals in 1968 and has continued to the present. The data set contains over 5,000 variables, including employment, income, and human-capital variables, as well as information on housing, travel to work, and mobility. In addition to the NLS and PSID data sets there are several other panel data sets that are of interest to economists, and these have been cataloged and discussed by Borus (1981) and Juster (2000); also see Ashenfelter and Solon (1982) and Beckett et al. (1988).¹

In Europe, various countries have their annual national or more frequent surveys – the Netherlands Socio-Economic Panel (SEP), the German Social Economics Panel (GSOEP), the Luxembourg Social Economic Panel (PSELL),

Table 1.1. *The NLS: Survey groups, sample sizes, interview years, and survey status*

Survey group	Age cohort	Birth year cohort	Original sample	Initial year/latest year	Number of surveys	Number at last interview	Status
Older men	45–59	4/2/07–4/1/21	5,020	1966/1990	13	2,092 ¹	Ended
Mature women	30–44	4/2/23–4/1/37	5,083	1967/1999	19	2,466 ²	Continuing
Young men	14–24	4/2/42–4/1/52	5,225	1966/1981	12	3,398	Ended
Young women	14–24	1944–1954	5,159	1968/1999	20	2,900 ²	Continuing
NLSY79	14–21	1957–1964	12,686 ³	1979/1998	18	8,399	Continuing
NLSY79 children	birth–14	—	— ⁴	1986/1998	7	4,924	Continuing
NLSY79 young adults	15–22	—	— ⁴	1994/1998	3	2,143	Continuing
NLSY97	12–16	1980–1984	8,984	1997/1999	3	8,386	Continuing

¹Interviews in 1990 were also conducted with 2,206 widows or other next-of-kin of deceased respondents.

²Preliminary numbers.

³After dropping the military (in 1985) and economically disadvantaged non-Black, non-Hispanic oversamples (in 1991), the sample contains 9,964 respondents eligible for interview.

⁴The sizes of the NLSY79 children and young adult samples are dependent on the number of children born to female NLSY79 respondents, which is increasing over time.

Source: NLS Handbook, 2000, U.S. Department of Labor, Bureau of Labor Statistics.

1.1 Advantages of Panel Data

3

the British Household Panel Survey (BHPS), etc. Starting in 1994, the National Data Collection Units (NDUs) of the Statistical Office of the European Communities, “in response to the increasing demand in the European Union for comparable information across the Member States on income, work and employment, poverty and social exclusion, housing, health, and many other diverse social indicators concerning living conditions of private households and persons” (Eurostat (1996)), have begun coordinating and linking existing national panels with centrally designed standardized multipurpose annual longitudinal surveys. For instance, the Mannheim Innovation Panel (MIP) and the Mannheim Innovation Panel – Service Sector (MIP-S) contain annual surveys of innovative activities (product innovations, expenditure on innovations, expenditure on R&D, factors hampering innovations, the stock of capital, wages, and skill structures of employees, etc.) of German firms with at least five employees in manufacturing and service sectors, started in 1993 and 1995, respectively. The survey methodology is closely related to the recommendations on innovation surveys manifested in the *OSLO Manual* of the OECD and Eurostat, thereby yielding international comparable data on innovation activities of German firms. The 1993 and 1997 surveys also become part of the European Community Innovation Surveys CIS I and CIS II (for detail, see Janz et al. (2001)). Similarly, the European Community Household Panel (ECHP) is to represent the population of the European Union (EU) at the household and individual levels. The ECHP contains information on demographics, labor-force behavior, income, health, education and training, housing, migration, etc. The ECHP now covers 14 of the 15 countries, the exception being Sweden (Peracchi (2000)). Detailed statistics from the ECHP are published in Eurostat’s reference data base New Cronos in three domains, namely health, housing, and income and living conditions (ILC).²

Panel data have also become increasingly available in developing countries. In these countries, there may not have a long tradition of statistical collection. It is of special importance to obtain original survey data to answer many significant and important questions. The World Bank has sponsored and helped to design many panel surveys. For instance, the Development Research Institute of the Research Center for Rural Development of the State Council of China, in collaboration with the World Bank, undertook an annual survey of 200 large Chinese township and village enterprises from 1984 to 1990 (Hsiao et al. (1998)).

Panel data sets for economic research possess several major advantages over conventional cross-sectional or time-series data sets (e.g., Hsiao (1985a, 1995, 2000)). Panel data usually give the researcher a large number of data points, increasing the degrees of freedom and reducing the collinearity among explanatory variables – hence improving the efficiency of econometric estimates. More importantly, longitudinal data allow a researcher to analyze a number of important economic questions that cannot be addressed using cross-sectional or time-series data sets. For instance, consider the following example taken from Ben-Porath (1973): Suppose that a cross-sectional sample of married women

4 Introduction

is found to have an average yearly labor-force participation rate of 50 percent. At one extreme this might be interpreted as implying that each woman in a homogeneous population has a 50 percent chance of being in the labor force in any given year, while at the other extreme it might imply that 50 percent of the women in a heterogeneous population always work and 50 percent never work. In the first case, each woman would be expected to spend half of her married life in the labor force and half out of the labor force, and job turnover would be expected to be frequent, with an average job duration of two years. In the second case, there is no turnover, and current information about work status is a perfect predictor of future work status. To discriminate between these two models, we need to utilize individual labor-force histories to estimate the probability of participation in different subintervals of the life cycle. This is possible only if we have sequential observations for a number of individuals.

The difficulties of making inferences about the dynamics of change from cross-sectional evidence are seen as well in other labor-market situations. Consider the impact of unionism on economic behavior (e.g., Freeman and Medoff 1981). Those economists who tend to interpret the observed differences between union and nonunion firms or employees as largely real believe that unions and the collective-bargaining process fundamentally alter key aspects of the employment relationship: compensation, internal and external mobility of labor, work rules, and environment. Those economists who regard union effects as largely illusory tend to posit that the real world is close enough to satisfying the conditions of perfect competition; they believe that the observed union–nonunion differences are mainly due to differences between union and nonunion firms or workers prior to unionism or postunion sorting. Unions do not raise wages in the long run, because firms react to higher wages (forced by the union) by hiring better-quality workers. If one believes the former view, the coefficient of the dummy variable for union status in a wage or earning equation is a measure of the effect of unionism. If one believes the latter view, then the dummy variable for union status could be simply acting as a proxy for worker quality. A single cross-sectional data set usually cannot provide a direct choice between these two hypotheses, because the estimates are likely to reflect interindividual differences inherent in comparisons of *different* people or firms. However, if panel data are used, one can distinguish these two hypotheses by studying the wage differential for a worker moving from a nonunion firm to a union firm, or vice versa. If one accepts the view that unions have no effect, then a worker's wage should not be affected when he moves from a nonunion firm to a union firm, if the quality of this worker is constant over time. On the other hand, if unions truly do raise wages, then, holding worker quality constant, the worker's wage should rise as he moves to a union firm from a nonunion firm. By following given individuals or firms over time as they change status (say from nonunion to union, or vice versa), one can construct a proper recursive structure to study the before–after effect.

1.1 Advantages of Panel Data

5

Whereas microdynamic and macrodynamic effects typically cannot be estimated using a cross-sectional data set, a single time-series data set usually cannot provide precise estimates of dynamic coefficients either. For instance, consider the estimation of a distributed-lag model:

$$y_t = \sum_{\tau=0}^h \beta_{\tau} x_{t-\tau} + u_t, \quad t = 1, \dots, T, \quad (1.1.1)$$

where x_t is an exogenous variable and u_t is a random disturbance term. In general, x_t is near x_{t-1} , and still nearer $2x_{t-1} - x_{t-2} = x_{t-1} + (x_{t-1} - x_{t-2})$; fairly strict multicollinearities appear among $h + 1$ explanatory variables, $x_1, x_{t-1}, \dots, x_{t-h}$. Hence, there is not sufficient information to obtain precise estimates of any of the lag coefficients without specifying, a priori, that each of them is a function of only a very small number of parameters [e.g., Almon lag, rational distributed lag (Malinvaud (1970))]. If panel data are available, we can utilize the interindividual differences in x values to reduce the problem of collinearity; this allows us to drop the ad hoc conventional approach of constraining the lag coefficients $\{\beta_{\tau}\}$ and to impose a different prior restriction to estimate an unconstrained distributed-lag model.

Another example is that measurement errors can lead to unidentification of a model in the usual circumstance. However, the availability of multiple observations for a given individual or at a given time may allow a researcher to identify an otherwise unidentified model (e.g., Biørn (1992); Griliches and Hausman (1986); Hsiao (1991b); Hsiao and Taylor (1991); Wansbeek and Koning (1989)).

Besides the advantage that panel data allow us to construct and test more complicated behavioral models than purely cross-sectional or time-series data, the use of panel data also provides a means of resolving or reducing the magnitude of a key econometric problem that often arises in empirical studies, namely, the often heard assertion that the real reason one finds (or does not find) certain effects is the presence of omitted (mismeasured or unobserved) variables that are correlated with explanatory variables. By utilizing information on both the intertemporal dynamics and the individuality of the entities being investigated, one is better able to control in a more natural way for the effects of missing or unobserved variables. For instance, consider a simple regression model:

$$y_{it} = \alpha^* + \boldsymbol{\beta}' \mathbf{x}_{it} + \boldsymbol{\rho}' \mathbf{z}_{it} + u_{it}, \quad \begin{array}{l} i = 1, \dots, N, \\ t = 1, \dots, T, \end{array} \quad (1.1.2)$$

where \mathbf{x}_{it} and \mathbf{z}_{it} are $k_1 \times 1$ and $k_2 \times 1$ vectors of exogenous variables; α^* , $\boldsymbol{\beta}$, and $\boldsymbol{\rho}$ are 1×1 , $k_1 \times 1$, and $k_2 \times 1$ vectors of constants respectively; and the error term u_{it} is independently, identically distributed over i and t , with mean zero and variance σ_u^2 . It is well known that the least-squares regression of y_{it} on \mathbf{x}_{it} and \mathbf{z}_{it} yields unbiased and consistent estimators of α^* , $\boldsymbol{\beta}$, and $\boldsymbol{\rho}$. Now suppose that \mathbf{z}_{it} values are unobservable, and the covariances between \mathbf{x}_{it} and \mathbf{z}_{it} are nonzero. Then the least-squares regression coefficients of y_{it} on

6 Introduction

\mathbf{x}_{it} are biased. However, if repeated observations for a group of individuals are available, they may allow us to get rid of the effect of \mathbf{z} . For example, if $\mathbf{z}_{it} = \mathbf{z}_i$ for all t (i.e., \mathbf{z} values stay constant through time for a given individual but vary across individuals), we can take the first difference of individual observations over time and obtain

$$y_{it} - y_{i,t-1} = \boldsymbol{\beta}'(\mathbf{x}_{it} - \mathbf{x}_{i,t-1}) + (u_{it} - u_{i,t-1}), \quad \begin{array}{l} i = 1, \dots, N, \\ t = 2, \dots, T. \end{array} \quad (1.1.3)$$

Similarly, if $\mathbf{z}_{it} = \mathbf{z}_t$ for all i (i.e., \mathbf{z} values stay constant across individuals at a given time, but exhibit variation through time), we can take the deviation from the mean across individuals at a given time and obtain

$$y_{it} - \bar{y}_t = \boldsymbol{\beta}'(\mathbf{x}_{it} - \bar{\mathbf{x}}_t) + (u_{it} - \bar{u}_t), \quad \begin{array}{l} i = 1, \dots, N, \\ t = 1, \dots, T, \end{array} \quad (1.1.4)$$

where $\bar{y}_t = (1/N) \sum_{i=1}^N y_{it}$, $\bar{\mathbf{x}}_t = (1/N) \sum_{i=1}^N \mathbf{x}_{it}$, and $\bar{u}_t = (1/N) \sum_{i=1}^N u_{it}$. Least-squares regression of (1.1.3) or (1.1.4) now provides unbiased and consistent estimates of $\boldsymbol{\beta}$. Nevertheless if we have only a single cross-sectional data set ($T = 1$) for the former case ($\mathbf{z}_{it} = \mathbf{z}_i$), or a single time-series data set ($N = 1$) for the latter case ($\mathbf{z}_{it} = \mathbf{z}_t$), such transformations cannot be performed. We cannot get consistent estimates of $\boldsymbol{\beta}$ unless there exist instruments that are correlated with \mathbf{x} but are uncorrelated with \mathbf{z} and u .

MaCurdy's (1981) work on the life-cycle labor supply of prime-age males under certainty is an example of this approach. Under certain simplifying assumptions, MaCurdy shows that a worker's labor-supply function can be written as (1.1.2), where y is the logarithm of hours worked, x is the logarithm of the real wage rate, and z is the logarithm of the worker's (unobserved) marginal utility of initial wealth, which, as a summary measure of a worker's lifetime wages and property income, is assumed to stay constant through time but to vary across individuals (i.e., $z_{it} = z_i$). Given the economic problem, not only is x_{it} correlated with z_i , but every economic variable that could act as an instrument for x_{it} (such as education) is also correlated with z_i . Thus, in general, it is not possible to estimate $\boldsymbol{\beta}$ consistently from a cross-sectional data set,³ but if panel data are available, one can consistently estimate $\boldsymbol{\beta}$ by first-differencing (1.1.2).

The "conditional convergence" of the growth rate is another example (e.g., Durlauf (2001); Temple (1999)). Given the role of transitional dynamics, it is widely agreed that growth regressions should control for the steady state level of income (e.g., Barro and Sala-i-Martin (1995); Mankiew, Romer, and Weil (1992)). Thus, a growth-rate regression model typically includes investment ratio, initial income, and measures of policy outcomes like school enrollment and the black-market exchange-rate premium as regressors. However, an important component, the initial level of a country's technical efficiency, z_{i0} , is omitted because this variable is unobserved. Since a country that is less efficient

1.1 Advantages of Panel Data

7

is also more likely to have lower investment rate or school enrollment, one can easily imagine that z_{i0} is correlated with the regressors and the resulting cross-sectional parameter estimates are subject to omitted-variable bias. However, with panel data one can eliminate the influence of initial efficiency by taking the first difference of individual country observations over time as in (1.1.3).

Panel data involve two dimensions: a cross-sectional dimension N , and a time-series dimension T . We would expect that the computation of panel data estimators would be more complicated than the analysis of cross-section data alone (where $T = 1$) or time series data alone (where $N = 1$). However, in certain cases the availability of panel data can actually simplify the computation and inference. For instance, consider a dynamic Tobit model of the form

$$y_{it}^* = \gamma y_{i,t-1}^* + \beta x_{it} + \epsilon_{it} \quad (1.1.5)$$

where y^* is unobservable, and what we observe is y , where $y_{it} = y_{it}^*$ if $y_{it}^* > 0$ and 0 otherwise. The conditional density of y_{it} given $y_{i,t-1} = 0$ is much more complicated than the case if $y_{i,t-1}^*$ is known, because the joint density of $(y_{it}, y_{i,t-1})$ involves the integration of $y_{i,t-1}^*$ from $-\infty$ to 0. Moreover, when there are a number of censored observations over time, the full implementation of the maximum likelihood principle is almost impossible. However, with panel data, the estimation of γ and β can be simplified considerably by simply focusing on the subset of data where $y_{i,t-1} > 0$, because the joint density of $f(y_{it}, y_{i,t-1})$ can be written as the product of the conditional density $f(y_{it} | y_{i,t-1})$ and the marginal density of $y_{i,t-1}$. But if $y_{i,t-1}^*$ is observable, the conditional density of y_{it} given $y_{i,t-1} = y_{i,t-1}^*$ is simply the density of ϵ_{it} (Arellano, Bover, and Labeager (1999)).

Another example is the time-series analysis of nonstationary data. The large-sample approximation of the distributions of the least-squares or maximum likelihood estimators when $T \rightarrow \infty$ are no longer normally distributed if the data are nonstationary (e.g., Dickey and Fuller (1979, 1981); Phillips and Durlauf (1986)). Hence, the behavior of the usual test statistics will often have to be inferred through computer simulations. But if panel data are available, and observations among cross-sectional units are independent, then one can invoke the central limit theorem across cross-sectional units to show that the limiting distributions of many estimators remain asymptotically normal and the Wald-type test statistics are asymptotically chi-square distributed (e.g., Binder, Hsiao, and Pesaran (2000); Levin and Lin (1993); Pesaran, Shin, and Smith (1999), Phillips and Moon (1999, 2000); Quah (1994)).

Panel data also provide the possibility of generating more accurate predictions for individual outcomes than time-series data alone. If individual behaviors are similar conditional on certain variables, panel data provide the possibility of learning an individual's behavior by observing the behavior of others, in addition to the information on that individual's behavior. Thus, a more accurate description of an individual's behavior can be obtained by pooling the data

8 Introduction

(e.g., Hsiao and Mountain (1994); Hsiao and Tahmiscioglu (1997); Hsiao et al. (1989); Hsiao, Applebe, and Dineen (1993)).

1.2 ISSUES INVOLVED IN UTILIZING PANEL DATA

1.2.1 Heterogeneity Bias

The oft-touted power of panel data derives from their theoretical ability to isolate the effects of specific actions, treatments, or more general policies. This theoretical ability is based on the assumption that economic data are generated from controlled experiments in which the outcomes are random variables with a probability distribution that is a smooth function of the various variables describing the conditions of the experiment. If the available data were in fact generated from simple controlled experiments, standard statistical methods could be applied. Unfortunately, most panel data come from the very complicated process of everyday economic life. In general, different individuals may be subject to the influences of different factors. In explaining individual behavior, one may extend the list of factors ad infinitum. It is neither feasible nor desirable to include all the factors affecting the outcome of all individuals in a model specification, since the purpose of modeling is not to mimic the reality but is to capture the essential forces affecting the outcome. It is typical to leave out those factors that are believed to have insignificant impacts or are peculiar to certain individuals.

However, when important factors peculiar to a given individual are left out, the typical assumption that economic variable y is generated by a parametric probability distribution function $P(y | \boldsymbol{\theta})$, where $\boldsymbol{\theta}$ is an m -dimensional real vector, identical for all individuals at all times, may not be a realistic one. Ignoring the individual or time-specific effects that exist among cross-sectional or time-series units but are not captured by the included explanatory variables can lead to parameter heterogeneity in the model specification. Ignoring such heterogeneity could lead to inconsistent or meaningless estimates of interesting parameters. For example, consider a simple model postulated as

$$y_{it} = \alpha_i^* + \beta_i x_{it} + u_{it}, \quad \begin{array}{l} i = 1, \dots, N, \\ t = 1, \dots, T, \end{array} \quad (1.2.1)$$

where x is a scalar exogenous variable ($k_1 = 1$) and u_{it} is the error term with mean zero and constant variance σ_u^2 . The parameters α_i^* and β_i may be different for different cross-sectional units, although they stay constant over time. Following this assumption, a variety of sampling distributions may occur. Such sampling distributions can seriously mislead the least-squares regression of y_{it} on x_{it} when all NT observations are used to estimate the model:

$$y_{it} = \alpha^* + \beta x_{it} + u_{it}, \quad \begin{array}{l} i = 1, \dots, N, \\ t = 1, \dots, T. \end{array} \quad (1.2.2)$$

1.2 Issues Involved in Utilizing Panel Data

9

For instance, consider the situation that the data are generated as either in case 1 or case 2.

Case 1: Heterogeneous intercepts ($\alpha_i^* \neq \alpha_j^*$), homogeneous slope ($\beta_i = \beta_j$). We use graphs to illustrate the likely biases due to the assumption that $\alpha_i^* \neq \alpha_j^*$ and $\beta_i = \beta_j$. In these graphs, the broken-line ellipses represent the point scatter for an individual over time, and the broken straight lines represent the individual regressions. Solid lines serve the same purpose for the least-squares regression of (1.2.2) using all NT observations. A variety of circumstances may arise in this case, as shown in Figures 1.1, 1.2, and 1.3. All of these figures depict situations in which biases arise in pooled least-squares estimates of (1.2.2) because of heterogeneous intercepts. Obviously, in these cases, pooled regression ignoring heterogeneous intercepts should never be used. Moreover, the direction of the bias of the pooled slope estimates cannot be identified a priori; it can go either way.

Case 2: Heterogeneous intercepts and slopes ($\alpha_i^* \neq \alpha_j^*$, $\beta_i \neq \beta_j$). In Figures 1.4 and 1.5 the point scatters are not shown, and the circled numbers signify the individuals whose regressions have been included in the analysis. For the example depicted in Figure 1.4, a straightforward pooling of all NT observations, assuming identical parameters for all cross-sectional units, would lead to nonsensical results because it would represent an average of coefficients that differ greatly across individuals. Nor does the case of Figure 1.5 make any sense in pooling, because it gives rise to the false inference that the pooled relation is curvilinear. In either case, the classic paradigm of the “representative agent” simply does not hold, and pooling the data under homogeneity assumption makes no sense.

These are some of the likely biases when parameter heterogeneities among cross-sectional units are ignored. Similar patterns of bias will also arise if the intercepts and slopes vary through time, even though for a given time period they are identical for all individuals. More elaborate patterns than those depicted here are, of course, likely to occur (e.g., Chesher and Lancaster 1983; Kuh 1963).

1.2.2 Selectivity Bias

Another frequently observed source of bias in both cross-sectional and panel data is that the sample may not be randomly drawn from the population. For example, the New Jersey negative income tax experiment excluded all families in the geographic areas of the experiment who had incomes above 1.5 times the officially defined poverty level. When the truncation is based on earnings, uses of the data that treat components of earnings (specifically, wages or hours) as dependent variables will often create what is commonly referred to as selection bias (e.g., Hausman and Wise (1977); Heckman (1976a, 1979); Hsiao (1974b)).

10 Introduction

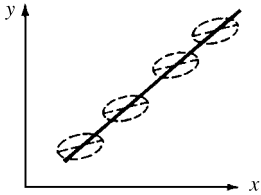


Fig. 1.1

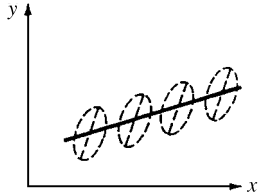


Fig. 1.2

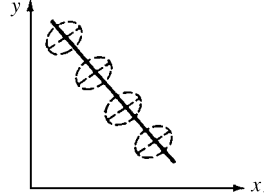


Fig. 1.3

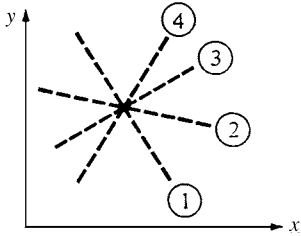


Fig. 1.4

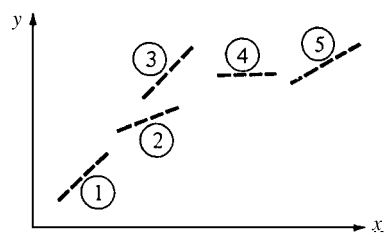


Fig. 1.5

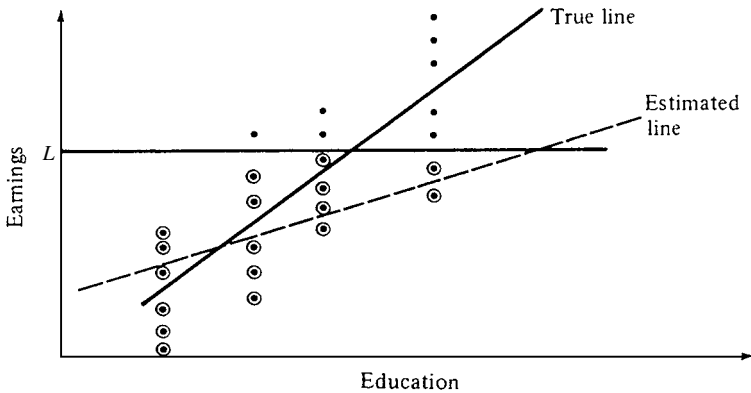


Fig. 1.6

For ease of exposition, we shall consider a cross-sectional example to get some idea of how using a nonrandom sample may bias the least-squares estimates. We assume that in the population the relationship between earnings (y) and exogenous variables (\mathbf{x}), including education, intelligence, and so forth, is of the form

$$y_i = \boldsymbol{\beta}'\mathbf{x}_i + u_i, \quad i = 1, \dots, N, \tag{1.2.3}$$

where the disturbance term u_i is independently distributed with mean zero and variance σ_u^2 . If the participants of an experiment are restricted to have earnings