

Andrews, Martyn J.; Schank, Thorsten; Upward, Richard

Working Paper

Practical estimation methods for linked employer-employee data

Diskussionspapiere, No. 29

Provided in Cooperation with:

Friedrich-Alexander-University Erlangen-Nuremberg, Chair of Labour and Regional Economics

Suggested Citation: Andrews, Martyn J.; Schank, Thorsten; Upward, Richard (2004) : Practical estimation methods for linked employer-employee data, Diskussionspapiere, No. 29, Friedrich-Alexander-Universität Erlangen-Nürnberg, Lehrstuhl für Arbeitsmarkt- und Regionalpolitik, Nürnberg

This Version is available at:

<https://hdl.handle.net/10419/28308>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.

FRIEDRICH-ALEXANDER-UNIVERSITÄT ERLANGEN-NÜRNBERG

Lehrstuhl für VWL, insbes. Arbeitsmarkt- und Regionalpolitik
Professor Dr. Claus Schnabel

**Diskussionspapiere
Discussion Papers**

No. 29

Practical estimation methods for linked employer-employee data

M.J. ANDREWS, T. SCHANK AND R. UPWARD

SEPTEMBER 2004

ISSN 1615-5831

PRACTICAL ESTIMATION METHODS FOR LINKED EMPLOYER-EMPLOYEE DATA^a

Martyn J. Andrews^b, Thorsten Schank^c, Richard Upward^d

ABSTRACT: Methods for the analysis of linked employer-employee data are not yet available in standard econometrics packages. In this paper, we make the fixed-effects methods developed originally by Abowd, Kramarz, Margolis and others more accessible, where possible, and show how they can be implemented in Stata. To illustrate these techniques, we give an example using German linked data. There is a caveat: when the number of plants is prohibitively large and the investigator wants to estimate the correlation between the worker and firm unobserved heterogeneities, the regression-based techniques discussed are not feasible. We also report an estimate of the correlation of zero.

ZUSAMMENFASSUNG: Die Analyse von zusammengeführten Personen- und Firmendaten ist bisher nicht in die Statistiksprogramme integriert worden. In dem vorliegenden Beitrag werden die ursprünglich von Abowd, Kramarz, Margolis u.a. entwickelten Analyseverfahren aufbereitet und, sofern möglich, wird gezeigt, wie diese in Stata implementiert werden können. Die vorgestellten Methoden werden mit einem kombinierten Firmen-Beschäftigtendatensatz (LIAB) aus Deutschland veranschaulicht. Es gibt jedoch eine Einschränkung: sofern die Anzahl der Firmen sehr groß ist und man die Korrelation zwischen den unbeobachtbaren Personen- und Firmenheterogenitäten schätzen möchte, können die in diesem Papier vorgestellten Regressionstechniken nicht verwendet werden. In Übereinstimmung mit anderen Studien finden wir ebenfalls eine Korrelation von Null.

KEYWORDS: linked employee-employer panel data, fixed effects

NEW JEL-CLASSIFICATION: C23, C87, J30

^aThe authors thank the IAB (Institut für Arbeitsmarkt- und Berufsforschung, Nürnberg) for kindly supplying the data, in particular, Lutz Bellmann and Stephan Bender. Financial support from the British Academy under Grant SG-35691 is also gratefully acknowledged. The views expressed in this paper are solely those of the authors and are not those of the IAB. Comments from presentations at the IAB, the Institute of Social and Economic Research at Essex, and the Departments of Economics at Manchester and Warwick are gratefully acknowledged. The usual disclaimer applies. All calculations were performed with Stata8/SE and the code given in the appendices is available from http://www.nottingham.ac.uk/economics/staff/details/richard_upward.html.

^bDr. Martyn J. Andrews, School of Economic Studies, University of Manchester, Manchester, M13 9PL, United Kingdom, martyn.andrews@man.ac.uk.

^cDr. Thorsten Schank, Friedrich-Alexander-Universität Erlangen-Nürnberg, Lehrstuhl für Arbeitsmarkt- und Regionalpolitik, Lange Gasse 20, D-90403 Nürnberg, thorsten.schank@wiso.uni-erlangen.de.

^dDr. Richard Upward, School of Economics, University of Nottingham, Nottingham, N67 2RD, United Kingdom, richard.upward@nottingham.ac.uk.

1 INTRODUCTION

Labour market outcomes are driven by the decisions of both workers and firms. However, it is only recently that the analysis of both sides of the market has become possible using matched (or linked) employer-employee data. There is a growing literature, whose origins are associated mainly with Abowd, Kramarz and Margolis. In Abowd, Kramarz & Margolis (1999) (hereafter AKM), they re-examine the whole of issue of persistent inter-industry wage differentials. Many other labour-market issues have been analysed, including inter-firm differences in productivity; the effects of hiring, quits and layoffs on productivity; the impact of new technology on wages; job creation and destruction; the effects of training; estimates of the cost of worker displacement; and the effects of unions and collective bargaining.²

Most econometric investigations of labour market issues are based on datasets that are either supply-side (individual- or household-level datasets) or demand-side (plant- or firm-level).³ If worker variables are correlated with firm variables, then any study that ignores information from the other side of the market will produce biased estimates. Biases also occur if the worker heterogeneity or the firm heterogeneity are correlated with the observables. For example, in AKM's (1999) paper 'High wage workers, high wage firms', it is unobservably better workers, in terms of wages, that are assumed to work in unobservably better firms.

Although there is a growing literature, the analysis of linked employer-employee data is not yet routine. There are two reasons why this research agenda has not moved on as quickly as it might. First, matched datasets involve linking together different sources of official information, and there are often technical, logistic and accessibility constraints that hinder progress. Second, there are various econometric issues to overcome, which mean that routine techniques and packages cannot be used. AKM's papers suggest these issues are quite technical. The objective of this paper, therefore, is to make these methods more accessible, where possible, and then show how they can be implemented in Stata. To illustrate these techniques, we give an example using German linked data, from the Institut für Arbeitsmarkt- und Berufsforschung, Nürnberg (hereafter IAB).⁴

A puzzle has emerged, in that the unobserved component of workers' wages appears to be *negatively* correlated with the unobserved component of firms' average wages. Apart from AKM's original study, which reported a positive correlation, all subsequent estimates have been negative. Abowd, Creecy & Kramarz (2002) (hereafter ACK) report that this is because the approximation used in their earlier work gives

²See also Abowd & Kramarz (1999) and Haltiwanger, Lane, Spletzer, Theeuwes & Troske (1999) for early surveys of the wide range of issues covered in this literature.

³Some datasets ask questions about the other side of the market; for example, a firm identifier and plant-size is available in the BHPS. Also, in what follows, 'workers' and 'individuals' are synonyms.

⁴Hereafter we refer to the data as LIAB: Linked IAB data.

different estimates when the models are re-estimated with the exact solution developed subsequently. Abowd, Creedy & Kramarz report correlations of -0.283 for French data and -0.025 for data from Washington State. Goux & Maurin (1999) find a correlation ranging from $+0.01$ to -0.32 depending on the time period chosen. Gruetter & Lalive (2003) find a correlation of -0.543 for Austrian data; Barth & Dale-Olsen (2003) report a correlation of between -0.47 and -0.55 . Our own estimates from German data suggest a correlation of approximately zero.

The paper is organised as follows. In Section 2, we set out the generic model that best represents the econometrics of fixed-effects models using matched employee-employer data, and in Section 3 we describe the various methods that can be used to estimate this generic model. In Section 4, we describe the LIAB data that we use to illustrate these techniques, which are presented in Sections 5 and 6. Section 7 concludes. Two appendices give the Stata code that can be used to estimate the models discussed in the paper.

2 A GENERIC MODEL

Consider the following model with both employer and employee heterogeneity and employer and employee covariates:

$$y_{it} = \mu + \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{w}_{jt}\boldsymbol{\gamma} + \mathbf{u}_i\boldsymbol{\eta} + \mathbf{q}_j\boldsymbol{\rho} + \alpha_i + \phi_j + \epsilon_{it} \quad (1)$$

There are $i = 1, \dots, N$ workers (N is often millions) and $j = 1, \dots, J$ firms (J is often thousands); y_{it} is the dependent variable; \mathbf{x}_{it} and \mathbf{u}_i are vectors of observable i -level covariates; \mathbf{w}_{jt} and \mathbf{q}_j are vectors of observable j -level covariates; and α_i and ϕ_j are (scalar) unobserved heterogeneities, correlated with observables and each other. Note that both α_i and \mathbf{u}_i are variables that are time-invariant for workers and similarly ϕ_j and \mathbf{q}_j are fixed over time for firms. \mathbf{x}_{it} , on the other hand, varies across i and t , and \mathbf{w}_{jt} varies across j and t . (There is more on use of j subscript below.) Equation (1) therefore contains all four possible types of information which a researcher might have about workers and firms. There are K observed covariates in total.

Both workers and firms are assumed to enter and exit the panel, which means we have an unbalanced panel with T_i observations per worker. There are $N^* = \sum_{i=1}^N T_i$ observations (worker-years) in total. Workers also change firms. This is crucial, as fixed-effects methods are identified by changers. In this paper, we assume ϵ_{it} is strictly exogenous, which implies that workers' mobility decisions are independent of ϵ_{it} . It is worth noting that mobility may be a function of the observables and the time-invariant unobservables.

Suppose the investigator only has access to worker (or household) data, and therefore considers estimating

$$y_{it} = \mu + \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{u}_i\boldsymbol{\eta} + \alpha_i + \phi_j + \epsilon_{it}.$$

If the investigator does not observe the vector $[\mathbf{w}_{jt}, \mathbf{q}_j]$ then the estimates of $[\boldsymbol{\beta}, \boldsymbol{\eta}]$ are biased if the vector $[\mathbf{x}_{it}, \mathbf{u}_i]$ is correlated with these missing firm-level variables. However, he can still control for ϕ_j providing he knows the identity of the firm, as there are multiple observations on workers within the same firm, which means that there are no biases arising from ϕ_j being correlated with any of the observables. Now suppose the investigator only has access to a single cross section. Clearly, he can still control for ϕ_j , but now he cannot control for α_i as it is now part of the error term $\alpha_i + \epsilon_i$.

Now suppose the investigator has only firm-level data, and considers estimating:

$$y_{jt} = \mu + \mathbf{w}_{jt}\boldsymbol{\gamma} + \mathbf{q}_j\boldsymbol{\rho} + \phi_j + \alpha_{jt} + \epsilon_{jt}.$$

Now the unit of observation is a firm, which means that $[y_{jt}, \alpha_{jt}, \epsilon_{jt}]$ are averages over each firm's employees. If everything were observed, *including* the vector of worker-level variables $[\mathbf{x}_{jt}, \mathbf{u}_j]$ (e.g. average age of the firm's employees, or the proportion of males in the firm), then the aggregation of variables would just cause heteroskedasticity. However, not observing $[\mathbf{x}_{jt}, \mathbf{u}_j]$ causes bias if these variables are correlated with the vector $[\mathbf{w}_{jt}, \mathbf{q}_j]$. However, we can control for ϕ_j using firm-level fixed effects methods, but we cannot control for α_{jt} , because it is part of the error term $\alpha_{jt} + \epsilon_{jt}$. This is the well-known aggregation bias caused by having firm-level rather than worker-level data.⁵ To conclude, without linked data, there are obvious biases from not observing observables, and from not controlling for unobservables.

Turning back to Equation (1), we emphasise that it is usual to assume that the heterogeneity terms α_i and ϕ_j are correlated with the observables. This means that random effects methods are inconsistent, and so fixed effects methods are needed to estimate the parameters of interest. This, in turn, means that $[\boldsymbol{\rho}, \boldsymbol{\eta}]$, the parameter vector associated with the time-invariant variables, is not identified. Rather than dropping $[\mathbf{u}_i, \mathbf{q}_j]$, it is usual to define

$$\theta_i \equiv \alpha_i + \mathbf{u}_i\boldsymbol{\eta} \tag{2}$$

and

$$\psi_j \equiv \phi_j + \mathbf{q}_j\boldsymbol{\rho} \tag{3}$$

giving

$$y_{it} = \mu + \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{w}_{jt}\boldsymbol{\gamma} + \theta_i + \psi_j + \epsilon_{it}. \tag{4}$$

⁵Early estimates of the union wage differential in the UK came from plant-level data (WIRS), which typically did not have important information on the employees' backgrounds.

Estimates of $[\boldsymbol{\eta}, \boldsymbol{\rho}]$ can be recovered by making the additional random effects assumptions $\text{Cov}(\mathbf{u}_i, \alpha_i) = \text{Cov}(\mathbf{q}_j, \phi_j) = 0$. Hausman & Taylor (1981) show that it is possible to identify time-varying effects using fixed-effects methods whilst identifying non-time-varying effects using random-effects methods in the same regression. However, some investigators may be unhappy about having different assumptions depending on whether the variable is time-invariant, or otherwise, so in everything that follows, we consider the identification of $[\boldsymbol{\eta}, \boldsymbol{\rho}]$ as an optional extra rather than part of the main story.

3 ECONOMETRIC METHODS

Equation (4) is the generic model that represents most of the existing literature. A number of fixed-effects methods have been proposed in the literature. In what follows, we describe each.⁶

3.1 LEAST SQUARES DUMMY VARIABLES (LSDV)

AKM are the first to propose consistent estimates of the parameters of Equations (1–4). It needs emphasising that they are particularly interested in estimating θ_i and ψ_j , in addition to $[\boldsymbol{\beta}, \boldsymbol{\gamma}]$, for two reasons. The first is that they want to see whether estimates of θ_i and ψ_j are correlated, hence the title ‘High wage workers, high wage firms’. The second is that they want to recover estimates of $\boldsymbol{\rho}$ and $\boldsymbol{\eta}$ using Equations (2) and (3) respectively. Because the heterogeneity variables are assumed to be correlated with the observables, they note that the Least Squares Dummy Variables (LSDV) estimator has the best properties, for the usual reasons. The LSDV estimates of α_i are inconsistent, although unbiased. (See Wooldridge (2002, ch. 10) for assumptions and properties of panel data models.) The properties of the ψ_j are the same as for $[\boldsymbol{\beta}, \boldsymbol{\gamma}]$, the parameters associated with the time-varying covariates $[\mathbf{x}_{it}, \mathbf{w}_{jt}]$.

There are two potential problems with actually computing this LSDV estimator. It is well known that a model with individual and time dummies (Baltagi’s Two Way Fixed Effects Model, Section 3.2) gives algebraic solutions for the estimates of the effects of the covariates *and* both sets of dummies. Essentially, there is a matrix that sweeps out both sets of dummies in one go, which means that a regression involving transformed variables is performed. For the model here, there are two important differences. First, in Baltagi the data are balanced, whereas here both workers and firms can enter and exit the panel. Wansbeek & Kapteyn (1989) analyse Baltagi’s model for unbalanced data, and obtain inelegant expressions that involve generalised inverses. Second, there is not a regular pattern between the firm and worker dummies

⁶The Stata code to estimate each of them can be downloaded from <http://www.arbeitsmarkt.wiso.uni-erlangen.de/schank.htm>.

as there is between Baltagi's individual and time dummies. It is the second that is the important difference, because it means that there is no algebraic transformation of the observables that sweeps away both heterogeneity terms in one go *and* which allows them to be recovered subsequently. To circumvent this second problem, AKM note that explicitly including dummy variables for the firm heterogeneity, but sweeping out the worker heterogeneity algebraically, gives exactly the same solution as the LSDV estimator.⁷

More precisely, the investigator must generate a dummy variable for each firm:

$$F_{it}^j = 1(j(i, t) = j) \quad j = 1, \dots, J,$$

where $1(\cdot)$ is the dummy variable indicator function and the function $j(i, t) = j$ maps worker i at time t to firm j . Now substitute

$$\psi_{j(it)} = \sum_{j=1}^J \psi_j F_{it}^j \quad (5)$$

into Equation (4).⁸ The θ_i are removed by time-demeaning (or differencing) over i :

$$y_{it} - \bar{y}_i = (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)\boldsymbol{\beta} + (\mathbf{w}_{jt} - \bar{\mathbf{w}}_i)\boldsymbol{\gamma} + \sum_{j=1}^J \psi_j (F_{it}^j - \bar{F}_i^j) + \epsilon_{it}, \quad (6)$$

where $\bar{z}_i = \sum_t z_{it}/T_i$ for any variable z . This means that J de-means (or differenced) firm dummies actually need creating.⁹ To distinguish this estimator from LSDV above, hereafter we label this estimator *FEiLSDVj*. They are identical estimators, but differ in how they are computed. The covariance matrix for FEiLSDVj needs the standard degrees-of-freedom adjustment, the formula for which is given in the next subsection.

We should note that $(F_{it}^j - \bar{F}_i^j)$ will be zero for all J dummies for any worker i who does not change firm. Furthermore, if we have a sample of firms, it will only be non-zero for workers who change from one firm within the sample to another firm in the sample. This means that for samples such as the LIAB, only a tiny proportion of workers have any non-zero terms. Identification of ψ_j is driven by the total number of such movers in each plant j . Some small plants may have no movers, in which case ψ_j is not identified. Other small plants may have only a very few movers, in which case estimates of ψ_j

⁷In linear models, there is no distinction between removing the heterogeneity algebraically or adding two full sets of dummy variables, for workers and firms, and so the terminology LSDV applies to both.

⁸Equation (5) shows that it would be better to use non-Greek letter for heterogeneity $\psi_{j(it)}$, because it is a variable, not a parameter.

⁹Differencing is ignored hereafter. There are various reasons why it is easier to implement the covariance transformation. Normally, the decision whether to estimate the model in first differences or use the covariance transform depends on which give the more efficient estimates. Both estimators are consistent. See Wooldridge (2002, Section 10.6.3).

will be very imprecise. This means that it may be not be sensible to estimate ψ_j for small firms, and instead one should group small firms together (this is what AKM and others do.)

To obtain estimates of the heterogeneity, first compute

$$\widehat{\psi}_{j(it)} = \sum_{j=1}^J \widehat{\psi}_j F_{it}^j \quad (7)$$

and then

$$\widehat{\theta}_i = \bar{y}_i - \overline{\widehat{\psi}}_i - \bar{\mathbf{x}}_i \widehat{\boldsymbol{\beta}} - \bar{\mathbf{w}}_i \widehat{\boldsymbol{\gamma}}$$

where $\overline{\widehat{\psi}}_i$ averages $\widehat{\psi}_{j(it)}$ over t .

There are two potential computational problems with this estimator. The first is the number of firms J , because the software needs to invert a matrix of dimension $(K + J) \times (K + J)$. For many applications, the number of firms is sufficiently small that FEiLSDVj is computationally feasible. For example, StataSE inverts 11,000 x 11,000 matrices. In our own empirical work, for reasons explained below, we only need to add approximately 2,000 firm dummies. There are many other situations where the number of firms/schools/doctors is sufficiently small. However, some datasets have tens of thousands of firms, or even hundreds of thousands (for, example, AKM and ACK). The second is the requirement that one must create and store J mean-deviations for N^* observations, meaning that the data matrix is $N^* \times (K + J)$. This may be prohibitively large for software packages which store all data in memory, such as Stata.

Some improvement in the storage efficiency of the J mean-deviated firm dummies can be achieved in Stata by using the lowest common multiple of all values of T_i . For example, if the data span a maximum of 5 years then T_i can be any value from $[1, 2, 3, 4, 5]$. Multiplying $F_{it}^j - \bar{F}_i^j$ by the lowest common multiple (in this case 60) yields a set of integers which can be stored in Stata as single bytes rather than 4- or 8-byte fractions.¹⁰

The memory requirements of the data matrix for the FEiLSDVj estimator are then approximately $(N^*J) + 4[N^*(K+1)]$ bytes. We require N^*J bytes for the mean-deviated firm dummies and $4[N^*(K+1)]$ bytes for the remaining K explanatory variables and the dependent variable, assuming each is stored as 4-bytes. In our example we have $N^* = 5,145,098$, $J = 1,821$ and $K = 64$, meaning that we require about 10GB of memory to proceed.

It is worth emphasising that firm dummies are no different from any multi-category dummy, so long as workers can move from category to another over time (e.g. region

¹⁰Storing the mean-deviated firm dummies as integers also appears to improve the accuracy of the matrix inversion procedure.

dummies, but not ethnicity dummies). This is why the notation \mathbf{w}_{jt} and \mathbf{q}_j is possibly confusing, since both are defined over every row indexed it . (Note that AKM use the notation $\mathbf{J}(i, t)$ to denote the mapping from worker i at time t to the firm j in which they are employed.) This means that the index j refers to the level of aggregation that w_{jt} actually varies over.

IDENTIFYING THE UNOBSERVED FIRM EFFECTS

An important issue is establishing how many unique unobserved firm effects can be identified. First, effects cannot be identified for firms which have no turnover; otherwise, $F_{it}^j - \bar{F}_i^j = 0$. Second, note that the firm dummies, when in mean-deviations, form a collinear set of variables

$$\sum_{j=1}^J (F_{it}^j - \bar{F}_i^j) = 0.$$

This is simply a consequence of having a collinear set of firm dummies, which sum to the constant before forming mean-deviations, and therefore sum to zero afterwards. In such a situation, one drops one of the firm dummies.

However, there is an additional identification issue, discussed by ACK. Identification of firm effects is only possible within a ‘group’, where a group is defined by the movement of workers between firms. A group contains all the workers who have ever worked for any of the firms in that group, and all the firms at which any of the workers were employed. A second (unconnected) group is defined only if no firm in the first group has ever employed any workers in the second, and no firms in the second group have ever employed any workers in the first. If there are G separate groups of firms, then it is not possible to identify one firm per group for the reason above.

ACK conclude that the number of estimable/identified person and firm effects is $N + J - G$, where N is the number of workers observed twice or more. Thus the correct degrees of freedom when estimating Equation (4) is $N^* - K - (J - G) - N$. When estimating Equation (6), the actual correct degrees of freedom are $N^* - K - (J - G)$, and so estimated standard errors need scaling by

$$\sqrt{\frac{N^* - K - (J - G)}{N^* - K - (J - G) - N}}. \quad (8)$$

A second implication of the grouping of firms is that estimates of $\hat{\psi}_j$ cannot be directly compared across groups. This is because it is arbitrary which ψ_j is set equal to zero for normalisation in each group. The same issue applies to the resulting $\hat{\theta}_i$. ACK suggest making the additional assumption that the average firm effect is the same across groups.

We have implemented the grouping algorithm in Stata.¹¹

IDENTIFYING THE EFFECTS OF TIME-INVARIANT VARIABLES

If the investigator can implement the LSDV estimator on i -de-meaned data (FEiLS-DVj), or implement one of AKM's other methods (discussed briefly below), AKM suggest that one can recover estimates of $\hat{\alpha}_i$ and $\hat{\phi}_j$ by estimating Equations (2, 3) as follows. First, run the auxiliary regressions:

$$\hat{\theta}_i = \text{const} + \mathbf{u}_i \boldsymbol{\eta} + \text{error} \quad (9)$$

$$\hat{\psi}_j = \text{const} + \mathbf{q}_j \boldsymbol{\rho} + \text{error} \quad (10)$$

which give consistent estimates of $\boldsymbol{\eta}$, $\boldsymbol{\rho}$ (AKM 1999, Section 3.4.4). Because α_i is dropped from (2), the identifying assumption is that $\text{Cov}(\mathbf{u}_i, \alpha_i) = 0$ or else there is omitted variable bias. Similarly, $\text{Cov}(\mathbf{q}_j, \phi_j) = 0$ is assumed in (3). One only needs N observations to estimate (2) and J observations to estimate (3). AKM estimate these equations by GLS, because of the aggregation to the firm-level. Because there are other causes of heteroskedasticity, one could use OLS and adjust the covariance matrix for clustering at the firm-level. Second, the investigator computes

$$\hat{\alpha}_i = \hat{\theta}_i - \mathbf{u}_i \hat{\boldsymbol{\eta}} \quad (11)$$

$$\hat{\phi}_j = \hat{\psi}_j - \mathbf{q}_j \hat{\boldsymbol{\rho}} \quad (12)$$

which are unbiased and asymptotic in T_i (Chamberlain 1984). θ and ψ can be defined at three levels of aggregation:

i, t	θ_i replicated T_i times	$\psi_{j(i,t)}$
i	θ_i	$\bar{\psi}_i = \sum_{t=1}^{T_i} \psi_{j(it)} / T_i$
j	$\bar{\theta}_j = \sum_{(it) \in j} \theta_i / N_j$	ψ_j

(N_j is the total number of worker-years observed in firm j .) AKM show that statistics based on aggregating $\hat{\theta}_i$ and $\hat{\alpha}_i$ to the level of the firm are consistent. To conclude, one can analyse distributions of $\hat{\psi}_j$, $\hat{\theta}_i$, specifically to see whether they are correlated.

¹¹The Stata code to calculate a grouping indicator for linked employer-employee data can be downloaded from <http://www.arbeitsmarkt.wiso.uni-erlangen.de/schank.htm>.

3.2 AKM'S APPROXIMATE METHODS

To deal with the large number of firm dummies, AKM propose a number of techniques in their (1999) paper that reduce the dimensionality of the problem. These require imposing further (testable) orthogonality assumptions. We do not discuss these further because ACK have recently developed a numerical solution for the LSDV estimator above.

3.3 ACK'S DIRECT LEAST SQUARES (DLS)

ACK, in addition to providing a more accessible discussion of their earlier papers, provide a numerical solution to the LSDV estimator of (1). They call it a Direct Least Squares Algorithm. They also make it clear that these methods are only relevant if one wants to estimate the heterogeneities. Finally, they re-estimate their original models on Washington and French data, and show that the AKM approximate methods reported in their (1999) paper give poorish results for the French data. Their solution involves an iterative technique that does not look easy to implement in standard software such as Stata.¹² More importantly, it is not regression based. The software is available from Abowd's website <http://instruct1.cit.cornell.edu/~jma7/abowdcv.html>.

3.4 SPELL FE

If one is not interested in the estimates of θ_i and ψ_j themselves, consistent estimates of β and γ from Equation (4) are straightforward to obtain by taking differences or by time-demeaning within each unique worker-firm combination (or 'spell'). This is because for each spell of a worker within a firm neither θ_i nor ψ_j vary. Defining $\lambda_s \equiv \theta_i + \psi_j$ as spell-level heterogeneity, which is swept out by subtracting averages at the spell-level, both θ_i and ψ_j have disappeared:

$$y_{it} - \bar{y}_s = (\mathbf{x}_{it} - \bar{\mathbf{x}}_s)\beta + (\mathbf{w}_{jt} - \bar{\mathbf{w}}_s)\gamma + (\epsilon_{it} - \bar{\epsilon}_s).$$

Again, the effects of \mathbf{u} and \mathbf{q} are not identified, because $\mathbf{u}_i - \bar{\mathbf{u}}_s = \mathbf{0}$ and $\mathbf{q}_j - \bar{\mathbf{q}}_s = \mathbf{0}$. In addition, any variable x_{it} or w_{jt} which is constant *within a spell* will also not be identified. One observation per spell is used up in identifying each spell fixed effect.¹³

This is basically the method that AKM discuss in Section 3.3, except they use differences rather than mean-deviations. AKM do not label this technique, so we call it

¹²Gruetter & Lalive (2003) also have an iterative technique, but it does not provide a covariance matrix.

¹³If there is just one observation per spell, then $y_{it} - \bar{y}_s = 0$, $\mathbf{x}_{it} - \bar{\mathbf{x}}_s = \mathbf{0}$, $\mathbf{w}_{jt} - \bar{\mathbf{w}}_s = \mathbf{0}$. This 'singleton' result can be used to reduce the sample size (by not much).

Spell FE or *FE(s)*. AKM state that it is consistent, inefficient, and “cannot be used to identify separately the firm intercept ... and the person effect”. It is clearly consistent as all the heterogeneity has been removed, and it is not the most efficient estimator because LSDV is. Because one cannot separate the worker and firm heterogeneities, AKM do not pursue this method further.

As when estimating any fixed-effects model, the standard errors may need correcting for the number of spells that the software has ‘forgotten’ about¹⁴

$$\sqrt{\frac{N^* - K}{N^* - K - S}}.$$

Unfortunately, given estimates of $\hat{\lambda}_s$, one cannot recover $\hat{\theta}_i$ and $\hat{\psi}_j$. Even if $S > N + J$, so that one could regress $\hat{\lambda}_s$ on worker and firm dummies, all that has happened is that β has been partitioned out of the problem, reducing the size of the problem by just K .

It is worth emphasising, however, that for many researchers this ‘spell fixed effects’ method is a practical and simple solution which does not present any computational difficulty, providing the investigator is not interested in analysing the heterogeneity post-estimation.

Spell FE is trivial to implement in Stata.¹⁵

IDENTIFYING THE EFFECTS OF TIME-INVARIANT VARIABLES: SPELL FEIV

We develop this method further to estimate the effects of time-constant variables \mathbf{u} and \mathbf{q} , which get swept away being constant within a spell. Consider the standard one-way fixed-effects model (say, using worker-level data only)

$$y_{it} = \mu + \mathbf{x}_{it}\beta + \theta_i + u_{it}. \quad (13)$$

The standard FE estimator of β can be interpreted as an IV estimator (Verbeek 2004, Section 10.2.5):

$$\begin{aligned} \hat{\beta}_{FE} &= [\Sigma_i \Sigma_t (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)' (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)]^{-1} \Sigma_i \Sigma_t (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)' (y_{it} - \bar{y}_i) \\ &= [\Sigma_i \Sigma_t (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)' \mathbf{x}_{it}]^{-1} \Sigma_i \Sigma_t (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)' y_{it} \end{aligned}$$

$x_{it} - \bar{x}_i$ is an ideal IV for any scalar x_{it} because: (i) it is uncorrelated with the unobservable θ_i , and (ii) it is correlated with x_{it} .

¹⁴Stata has a command `areg` which does not need this correction. Also, it can correct the standard errors for clustering, which, in this context, should be at the firm level. Wooldridge (2002, p. 57) explains why this correction is needed even when there is clustering.

¹⁵See do-file downloadable from <http://www.arbeitsmarkt.wiso.uni-erlangen.de/schank.htm>.

This implies one can estimate Equation (13) by IV GLS with $\mathbf{x}_{it} - \bar{\mathbf{x}}_i$ as an IV for \mathbf{x}_{it} . The other extreme case uses \mathbf{x}_{it} as an IV, which generates the random effects estimator. The objective here is to estimate the parameters of Equation (1), not Equation (4). The above argument implies that it is possible to estimate the parameters on the time-varying variables by time-demeaning them, *and* to estimate the parameters of the time-invariant variables using random effects. This approach can be thought of as ‘in between’ the FE estimator (which cannot estimate the parameters on time-invariant variables) and the RE estimator (which does not allow for any correlation between the time-varying variables and the unobservable heterogeneity). All variables that are correlated with unobservables (\mathbf{x}_{it} , \mathbf{w}_{jt}) are instrumented by their mean deviations $\mathbf{x}_{it} - \bar{\mathbf{x}}_s$ and $\mathbf{w}_{jt} - \bar{\mathbf{w}}_s$ respectively. This is not possible for the time invariant variables, (\mathbf{u} , \mathbf{q}), which can only be instrumented by themselves, which means we are assuming that $\text{Cov}(\mathbf{u}_i, \alpha_i) = 0$ and $\text{Cov}(\mathbf{q}_j, \phi_j) = 0$. In other words we are making exactly the same assumptions for \mathbf{u} and \mathbf{q} as we have done throughout, which is why Spell FEIV is a side-issue. This is a special case of Hausman & Taylor’s (1981) estimator.

3.5 TWO-STEP METHOD

The main problem with the FEiLSDVj estimator is that it requires the inversion of a $(K + J) \times (K + J)$ cross-product matrix. As noted, in some cases J may be only a few thousand, and so the estimator is feasible. This is particularly true where we have a sample of plants, and if we only attempt to identify the firm effects for larger firms. There is another constraint however, which is the sheer number of observations, even when J is sufficiently small. This is because the data matrix is $N^* \times (K + J)$, and might be prohibitively large for software packages that store data in memory rather than on disk. To circumvent this problem, we propose the following two-step method, based on the fact that only movers between firms identify firm effects.

In the first step, the investigator estimates Equation (6), but only using those observations that identify the ψ_j s. These are workers who move between firms. As above, compute $\hat{\psi}_{j(it)}$ (and $\hat{\theta}_i$) where they are identified. $\hat{\psi}_{j(it)}$ does not exist for plants that have no movers, and a normalisation restriction is needed for each ‘group’ of connected plants.

In the second step, the investigator estimates the following version of Equation (4), using *all* the data where $\hat{\psi}_{j(it)}$ exists:

$$y_{it} = \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{w}_{jt}\boldsymbol{\gamma} + \theta_i + \delta\hat{\psi}_{j(it)} + \epsilon_{it}.$$

In other words, the firm-level heterogeneity term ψ_j is replaced by estimates from the first stage (where δ is a scalar.) He then takes deviations from worker means, as per

usual:

$$y_{it} - \bar{y}_i = (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)\boldsymbol{\beta} + (\mathbf{w}_{jt} - \bar{\mathbf{w}}_i)\boldsymbol{\gamma} + \delta(\widehat{\psi}_{j(it)} - \overline{\widehat{\psi}}_i) + \epsilon_{it} \quad (14)$$

Using $\widehat{\psi}_{j(it)}$ and the second-step estimates, denoted $\widetilde{}$, the investigator then computes

$$\widetilde{\theta}_i = \bar{y}_i - \bar{\mathbf{x}}_i\widetilde{\boldsymbol{\beta}} - \bar{\mathbf{w}}_i\widetilde{\boldsymbol{\gamma}} - \overline{\widehat{\psi}}_i$$

$\widetilde{\theta}_i$ and $\widehat{\psi}_j$ can be analysed in the usual way, namely to compute the correlation between them, and to run regressions (9,10) above to recover estimates on the time-invariant variables.

3.6 A ROAD MAP?

To conclude the discussion of the methods discussed in this section, we outline a flow chart that should help the investigator decide which method is appropriate for his needs.

1. Does the investigator want to estimate employer and employee heterogeneity?

No Use Spell-level FE

Yes ...

2. Are there too many firm dummies to add ‘by hand’?

Yes Use AKM techniques

No ...

3. Is there enough memory?

Yes Use (transformed) firm dummies (FEiLSDVj)

No Use the Two-step method

For all methods, one can recover estimates on \mathbf{u}_i and \mathbf{q}_j making standard RE assumptions. The Stata code for estimating all of the models outlined in this section, apart from ACK’s DLS, can be downloaded from

<http://www.arbeitsmarkt.wiso.uni-erlangen.de/schank.htm>.

4 THE DATA (LIAB)

4.1 THE IAB ESTABLISHMENT PANEL (BETRIEBSPANEL)

The IAB (Institut für Arbeitsmarkt- und Berufsforschung) collect their own demand side data: the Betriebspanel is an establishment panel of $\approx 8,000$ establishments located in the former West Germany and $\approx 8,000$ establishments in the former East Germany. It covers the period 1993–present (1996–present for East Germany) and covers 1% of all plants and 7% of all employees in the population. The establishments are selected using a fairly complicated weighting procedure. (See Kölling (2000) for full details on the Betriebspanel.) Information on each establishment includes:¹⁶

- Total employment (also disagg'd) (`size1-size10`)
- Standard (`lhbar`) and overtime hours
- Wage recognition (`B,B1,B2`)
- Output
- Exports
- Investment (`inv`)
- Wage bill
- Urbanicity (`urban1-urban10`)
- Geographical location
- Nationality of ownership (foreign in 2000)
- Technology (subjective measure)
- Organisational change (subjective measure)
- Profitability (`profit1-profit5`)
- Age of plant (`vin`) and whether parent is a single-plant firm (`single`)

¹⁶If variables are used in tables below, their acronyms are also given. For full definitions, see Table 2.

4.2 THE EMPLOYMENT STATISTICS REGISTER (BESCHÄFTIGTENSTATISTIK)

On the other side of the labour market, the IAB has access to the employment statistics register (Beschäftigtenstatistik). It is an administrative panel of all employees who are covered by the social security system (about 80% of total employment), and is collected by the plant. There is at least one compulsory notification during each calendar year. It covers 1975–present for West Germany and 1992–present for East Germany. It contains about 400 million records, covering about 46 million employees. (See Bender, Haas & Klose (2000) for full details on the Beschäftigtenstatistik.) Information on each worker includes:

- Gender (`female`), age (`age`), nationality (`foreign`), marital status (`married`)
- Start and end dates of every employment spell (`mjob` for more than one job)
- Occupation (3-digit) (`occ1-occ6`)
- Daily wages (left truncated and right censored) (`lw`, but see below for more information)
- Qualifications: education/apprenticeship (`qual1-qual6`)
- Industry (`ind1-ind10`)
- Region
- Establishment identification number

4.3 THE LINKED IAB EMPLOYER-EMPLOYEE DATA (LIAB)

By using the establishment identification number, the IAB are able to associate each worker in the Beschäftigtenstatistik with an establishment in the IAB panel. Note that it is also possible to aggregate up all workers (not just those employed by establishments in the panel) to the establishment level. The particular dataset we use for this study was created by selecting all employees in the employment register who are employed by the surveyed establishments on June 30th each year. The data we use cover 1990–97 and contain 118,399,405 observations (*it* rows).

4.4 SAMPLE USED FOR WAGE EQUATIONS

To illustrate the techniques outlined above, we estimate various standard wage equations. The sample we use covers 1993–97, that is $1 \leq T_i \leq 5$, and is for West Germany only. We also drop observations for apprentices, part-timers, homeworkers and those

with a daily wage of less than 10 *DM*. In addition, the data are right-censored.¹⁷ As always, we also drop observations with missing values.

Workers change plants, and in particular, can change between plants that are surveyed in the IAB panel and plants that are not. In this study, we keep only those years (*it* rows) when a worker is working in an IAB-panel plant. This is because we don't observe \mathbf{w}_{jt} or \mathbf{q}_j in those years when a worker is working for a non-IAB plant. Table 1 summarises the data, in exactly the same format used by AKM.

[TABLE 1 ABOUT HERE]

Identification of unobserved plant-effects is driven only by those workers who change plants. Thus an important sub-sample comprises those workers who have two or more spells ($S_i > 1$) in IAB plants ('IAB movers'). In Table 1, workers who return to the same employer after an intervening spell with another employer are coded as starting a new spell. In Section 3.4, a spell is defined as any unique worker/employer combination, and so all periods a worker spends with a given employer are coded as a single spell. This is why there are 1,954,242 spells in Table 1 but only 1,953,774 spells in the regression sample. This, and the sample of IAB movers, is summarised below.

	all	IAB movers
No of obs	5,145,098	72,253
No of inds	1,930,260	23,393
No plants	4,376	1,821
Obs/inds	2.67	2.69
Inds/plant	441	997
No spells	1,953,774	46,635

The 72,253 observations comprise 23,272 workers with 2 spells and 121 workers with 3 spells. This makes 23,393 workers who have at least two spells. The total number of spells is $2 \times 23,272 + 3 \times 121 = 46,635$ spells, although only $23,272 + 2 \times 121 = 23,514$ spells are usable in spell-level fixed-effects regressions, as the first spell for each workers is not used (which is why workers with only one spell are not in this sample.)

Are these samples representative? As already discussed, the IAB-panel plants over-represent large plants in the population, and so workers in IAB plants are not a random sub-sample of the population. It is also possible that the 23,393 workers who move between IAB plants may not a random sub-sample of 1,930,260; exactly the same issue arises in all panel data models, which rely on movers for identification. (For example, estimates of union wage differential based on a sample of joiners/quitters.)

Table 2 reports sample means: the first three columns average by workers whereas columns four to six average by plants. For example, in column one, the regression

¹⁷In a paper that is concerned with methods, this is not as issue, although one could deal with this in the same way as Gruetter & Lalive (2003, p.6).

sample, 22.95% of workers are female whereas, on average, each plant employs 34.76% females (column four). These sets of means are often different from each other because of the underlying nature of the plant-size distribution. Workers are much more likely to work for large plants rather than small plants. Because large plants have higher wages, average log earnings are much smaller in column four than in column one. There are also big differences in sample means for whether married, qualifications, industry, union bargaining, investment and the age of the plant.

Column two corresponds to column one, but for the 23,982 workers who move between IAB plants. The difference between columns one and two is in column three. Column five corresponds to column four, but for the 1,821 plants that experience ‘IAB turnover’, that is employ workers who move between IAB plants. The difference between columns four and five is in column six. As we only identify 1,821 plants out of 4,376, the obvious question is whether these plants are *observably* the same, on average, as the 4,376? The same question applies to whether the 23,393 movers are observably the same as the 1,930,260 workers. In fact, the 1,821 plants pay lots more (0.1678 log-points), employ fewer females, employ more married workers, tend to be bigger firms located in different industries, and invest more (column six). Looking at individual workers, movers only get slightly more pay (0.0327 log-points), are younger, are less likely to be women, are more highly qualified, and are employed at plants with lower investment (column three).

[TABLE 2 ABOUT HERE]

Even if this sub-sample is not random, it does not follow that the estimates of $1,821 \hat{\psi}_j$ are inconsistent. This depends on what causes movement. If based on match quality, say $f(\alpha, \phi)$, then estimates are consistent because α, ϕ are swept away. However, it is a strong assumption to suggest that movement is independent of ϵ ; any shock that affects workers and firms suggests that movement and ϵ are correlated.

We conclude this discussion on the identification of unobserved plant-effects by counting the number of movers for each plant. Figure 1 plots the cumulative frequency for the number of plants against the number of movers. For example, one plant has 1,886 movers, but 472 plants only have one mover, and 2,555 plants have no movers at all. This is a very skewed distribution, and is a feature of linked employee-employer datasets. The IAB panel is a 1% sample of plants. Even though it is a large sample, the probability of observing a worker moving from one IAB plant to another is very small. Even if one observed the population of plants, very small plants would experience little or no turnover in a five-year period, making estimation of their ψ_j very noisy.

One possible strategy the investigator might adopt is to only identify ψ_j for plants with more than x movers, and group all remaining small plants into one plant (Abowd

et al. 2002). Using Figure 1, we set $x = 30$, giving 211 large plants and one small plant (albeit with a lot of employees). To conclude, it is important for the investigator to be aware of how little information is sometimes used to identify each unobserved plant effect, especially if plants are small.

5 RESULTS

[TABLE 3 ABOUT HERE]

Table 3 reports three conventional models, so described because they control for heterogeneity from only one side of the market, at best. The first is labelled Pooled OLS, which is Equation (1) where neither α_i nor ϕ_j are controlled for, of which there are three variants. The first only includes worker-level covariates, the second only plant-level covariates, and third includes both sets. The idea here is to assess the extent to which estimates on worker-level covariates are affected by the absence of plant-level covariates, and *vice versa*—in other words, to assess the extent to which the two sets of covariates are correlated with each other. A comparison of the estimates shows that the estimates do change, but not by much. The plant-level covariates move more, which is expected, given their standard errors are generally bigger.

The second model is labelled FE(i), i.e. the worker-level heterogeneity θ_i is controlled for, but ϕ_j becomes part of the model's error term:

$$y_{it} = \mu + \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{w}_{jt}\boldsymbol{\gamma} + \mathbf{q}_j\boldsymbol{\rho} + \theta_i + (\phi_j + \epsilon_{it})$$

Notice that the effects of the time-invariant worker-level variables \mathbf{u}_i are not identified, namely `foreign` and `female`. The extent to which an estimate moves compared with Pooled OLS (previous column) depends on the extent to which θ_i is correlated with observed covariates. Here there are some large movements. Notice that Stata reports an estimate of the correlation between the deterministic part of the regression and θ_i (`'corr(ui,Xb)'`), and there is very strong negative correlation of -0.66 , which is a different manifestation of the same thing.¹⁸

The third model is labelled FE(j)

$$y_{it} = \mu + \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{w}_{jt}\boldsymbol{\gamma} + \mathbf{u}_i\boldsymbol{\eta} + \psi_j + (\alpha_i + \epsilon_{it})$$

Now the plant-level heterogeneity ψ_j is controlled for, but α_i becomes part of the model's error term. The effects of the time-invariant plant-level variables \mathbf{q}_j are not identified, namely industry dummies and a dummy for whether the plant is `single`.

¹⁸`'corr(ui,Xb)'` varies from model to model. For FE(i), it is the correlation between $\hat{\theta}_i$ and $\mathbf{x}_{it}\hat{\boldsymbol{\beta}} + \mathbf{w}_{jt}\hat{\boldsymbol{\gamma}} + \mathbf{q}_j\hat{\boldsymbol{\rho}}$.

This is not a model that one would normally estimate, but is useful if the investigator *cannot* control for both ψ_j and α_i simultaneously, because it at least indicates the extent to which ψ_j is correlated with the observed covariates. Here the correlation between ψ_j and the deterministic part of the model is much weaker, and positive, at 0.08.

What is missing, of course, is that we do not control for any correlation of *both* unobserved fixed effects, ϕ_j and α_i , with observable characteristics. Table 4 reports two models that do exactly this.

[TABLE 4 ABOUT HERE]

The first of these is FE(s), the easy-to-use technique that removes spell-level heterogeneity (Section 3.4 above). The effects of all time-invariant covariates are not identified, but the estimates of the time-varying covariates are consistent. If the investigator is not interested in estimating the worker- and plant-heterogeneities, he can stop here. Comparing these estimates with Pooled OLS and FE(i) in the previous table is of some considerable interest, as these are the better estimates. Notice that the correlation between the deterministic part of the regression and λ_s is -0.56 , which is approximately equal to the sum of those from FE(i) and FE(j). Given that $\lambda_s = \theta_i + \psi_j$, this is not surprising. The IV version that estimates the effects of time-invariant variables, under the extra assumptions $\text{Cov}(\mathbf{u}_i, \alpha_i) = \text{Cov}(\mathbf{q}_j, \phi_j) = 0$, is reported in the second column. The estimates of the time-varying covariates are virtually identical.

Following the ‘road-map’ outlined in Section 3.6, the next decision is to ascertain whether there are too many plant dummies to add ‘by hand’ when estimating Equation (6). This technique, if feasible, is labelled FEiLSDVj. ‘By hand’ means that dummies for each plant are explicitly added to the regression like any other covariate; that is, cannot be dealt with algebraically. In the models being estimated here, we have 5,145,098 observations, and need $J - G = 1,821 - 33 = 1,788$ plant dummies, these being those plants which have IAB turnover, i.e. movers to/from another IAB plant. The memory needed is too prohibitive. As discussed on Page 8, we consider implementing the trick whereby we multiply the dummies by 60 so that they are stored as single bytes. This didn’t work: we have $N^* = 5,145,098$, $J = 1,821$ and $K = 64$, meaning that we still require about 10GB of memory to proceed.

Thus the only way forward is to use the Two-step method outlined in Section 3.5. This is the second model in Table 4. In the first step, we estimate Equation (6), but only using the 72,253 observations (23,393 IAB movers) that identify the 1,821 ψ_j s. This is reported in the column labelled ‘1st Step’. The estimates on the plant dummies are then used to form the single variable $\hat{\psi}_j$, using Equation (7). Note that we only identify 1,788 plants, not the 1,821 in the dataset. This is because there are 33 ‘groups’ of unconnected plants, and one plant per group is not identified. Each

$\widehat{\psi}_j$ is normalised on the average $\widehat{\psi}_j$ for its group g . In the second step, $\widehat{\psi}_j$ is added to a standard model with worker-level heterogeneity, which is removed using worker mean-deviations, i.e. Equation (14). This is reported in the column labelled ‘2nd Step’. Because 33 $\widehat{\psi}_j$ s are not identified, the resulting number of observations falls to 4,873,901, corresponding to 1,812,562 workers. The last two columns in Table 4 report estimates of the auxiliary regressions shown in Equations (9, 10), whereby estimates of the time-invariant covariates $[\mathbf{u}_i, \mathbf{q}_j]$ are recovered, under the usual assumptions $\text{Cov}(\mathbf{u}_i, \alpha_i) = \text{Cov}(\mathbf{q}_j, \phi_j) = 0$.

In the Two-step method, one would expect that the estimate on $\widehat{\psi}_{j(i,t)}$ to be close to unity, which might be viewed as some form of specification check. The actual estimate is 0.943, which is about three standard errors lower than unity.

The FE(s) and Two-step methods give very similar estimates of the time-varying covariates, which illustrates that the Two-step method also gives consistent estimates of β and γ . However, the estimates of the time-invariant covariates do differ, probably because the estimates $\widehat{\psi}_{j(i,t)}$ used as a dependent variable in the last column are unreliable, given the discussion on their identification above. The estimates for the $\widehat{\theta}_i$ regression are much closer to Spell FEIV.

As emphasised repeatedly, the advantage of the Two-Step method over FE(s) is that estimates of θ_i and ψ_j are obtained. The means of these two distributions are not identified, but estimates of their variances are easily computed, as is the correlation between them. It is the correlation that is particularly interesting, since it estimates the extent to which unobservably ‘good’ workers are employed in unobservably ‘good’ plants. The correlations between $\widehat{\psi}_j$, the first-step estimates $\widehat{\theta}_i$, and the second-step estimates $\widetilde{\theta}_i$, are reported below. Correlations with the corresponding estimates $\widetilde{\alpha}_i = \widetilde{\theta}_i - \mathbf{u}_i \widetilde{\eta}$ and $\widehat{\phi}_j = \widehat{\psi}_j - \mathbf{q}_j \widehat{\rho}$ are reported for completeness (see Equations (11,12)) but are of less interest:

	$\widetilde{\theta}$	$\widehat{\theta}$	$\widehat{\psi}$	$\widetilde{\alpha}$	$\widehat{\phi}$
$\widetilde{\theta}$	1.0000				
$\widehat{\theta}$	0.9291	1.0000			
$\widehat{\psi}$	-0.1684	-0.1608	1.0000		
$\widetilde{\alpha}$	0.9596	0.9055	-0.2097	1.0000	
$\widehat{\phi}$	-0.2021	-0.1846	0.9313	-0.2371	1.0000
Uses 4,883,331 <i>it</i> observations.					

The important finding is that $\text{corr}(\widehat{\psi}, \widetilde{\theta}) = -0.1684$. This correlation has the wrong sign if one expects that unobservably ‘good’ workers would be employed in unobservably ‘good’ plants. However, all of the literature (summarised briefly in the Introduction) finds a negative correlation, which gives rise to the question as to whether this a genuine economic phenomenon or whether there is a technical issue insofar as this

estimate is downwards biased. Our own view is that it is the latter (Andrews, Schank & Upward 2004), and that the size of the bias decreases with the number of movers used in estimating each $\hat{\psi}_j$.

Under the assumptions of the model, we have now consistent estimates of all the components of the RHS of Equation (4)

$$y_{it} = \mathbf{x}_{it}\hat{\boldsymbol{\beta}} + \mathbf{w}_{jt}\hat{\boldsymbol{\gamma}} + \hat{\theta}_i + \hat{\psi}_j + \hat{\epsilon}_{it}$$

where the hat now refers to any consistent estimate (Two-step, FEiLSDVj, or AKM's DLS). This allows us to analyse the correlations between the observed and unobserved components of wages, on both sides of the market:

	$\hat{\theta}$	$\hat{\psi}$	$\mathbf{x}_{it}\hat{\boldsymbol{\beta}}$	$\mathbf{w}_{jt}\hat{\boldsymbol{\gamma}}$
$\hat{\theta}$	1.0000			
$\hat{\psi}$	-0.1684	1.0000		
$\mathbf{x}_{it}\hat{\boldsymbol{\beta}}$	0.0235	0.0361	1.0000	
$\mathbf{w}_{jt}\hat{\boldsymbol{\gamma}}$	0.0306	-0.3198	0.0033	1.0000
Uses 4,883,331 <i>it</i> observations.				
	$\hat{\theta}$	$\hat{\psi}$	$\mathbf{x}_{it}\hat{\boldsymbol{\beta}}$	$\mathbf{w}_{jt}\hat{\boldsymbol{\gamma}}$
$\hat{\theta}$	1.0000			
$\hat{\psi}$	-0.1574	1.0000		
$\mathbf{x}_{it}\hat{\boldsymbol{\beta}}$	0.0148	0.0531	1.0000	
$\mathbf{w}_{jt}\hat{\boldsymbol{\gamma}}$	0.0028	-0.3558	-0.1138	1.0000
Averages to 1,816,368 <i>i</i> observations.				
	$\hat{\theta}$	$\hat{\psi}$	$\mathbf{x}_{it}\hat{\boldsymbol{\beta}}$	$\mathbf{w}_{jt}\hat{\boldsymbol{\gamma}}$
$\hat{\theta}$	1.0000			
$\hat{\psi}$	-0.4339	1.0000		
$\mathbf{x}_{it}\hat{\boldsymbol{\beta}}$	0.0508	0.1148	1.0000	
$\mathbf{w}_{jt}\hat{\boldsymbol{\gamma}}$	0.1006	-0.3974	-0.0267	1.0000
Averages to 1,821 <i>j</i> observations.				

Even though aggregating information to the plant-level means that estimators remain consistent, it is noticeable that correlations get bigger in absolute size. Looking at the *it*-level correlations, they generally make sense, except for those involving ψ . In particular, $\text{corr}(\psi, \mathbf{w}\hat{\boldsymbol{\gamma}}) = -0.3198$ looks somewhat awry, as well as $\text{corr}(\hat{\psi}, \hat{\theta})$ discussed above. The observed components are uncorrelated with each other, $\text{corr}(\mathbf{x}\hat{\boldsymbol{\beta}}, \mathbf{w}\hat{\boldsymbol{\gamma}}) = 0.0033$, which means that ignoring information from one side of the market does not affect estimates from the other side. All of the other cross-market correlations are small: $\text{corr}(\hat{\theta}, \mathbf{w}\hat{\boldsymbol{\gamma}}) = 0.0306$ and $\text{corr}(\hat{\psi}, \mathbf{x}\hat{\boldsymbol{\beta}}) = 0.0361$. Also, the unobserved and observed components of workers' wages are uncorrelated, $\text{corr}(\hat{\theta}, \mathbf{x}\hat{\boldsymbol{\beta}}) = 0.0235$. In short, it is the three correlations that involve $\hat{\psi}$ that looks wrong, and confirms these estimates of ψ are often 'poor', being identified from plants that have very little turnover.

To investigate this further, we group all but the smallest 211 plants into one plant, reported in Table 5.¹⁹ Now all the plants are connected, i.e. $G = 1$. One advantage of doing this is that we are able to estimate the model using FEiLSDVj (first three columns), which is the estimator with the best properties. This allows us to make two comparisons. The first is to re-estimate the model using the Two-Step method (final three columns), thereby compare the two estimation methods directly. The second is that this Two-Step method for a model with 212 plants can be compared with the same method applied to the model that has 1,821 plants, discussed immediately above and reported in Table 4. (Note that we do not report the 1st Step estimates in Table 5.)

The estimates reported in columns one and four are very similar to each other, as are the standard errors, which illustrates clearly that our Two-Step method has a lot to recommend it. One obvious reason why there are differences between the two methods is that the sample sizes differ, for reasons already explained. Also notice that the estimate on $\hat{\psi}_j$ is 0.9872, and is now insignificantly different from unity. The one place where the two models differ is that the estimates of the \mathbf{q}_j variables are somewhat bigger, which suggests that the estimates $\hat{\psi}_j$ for large plants exhibit sampling variation, and therefore differ across the two estimation methods. The sample correlation between the $\hat{\psi}_j$ s for both models is 0.459 (but is much bigger for the second-step estimates of θ).

The correlation between θ_i and ψ_j are -0.0172 for FEiLSDVj method and -0.0239 for the Two Step method. They are not exactly the same because the estimates of $\hat{\psi}_j$ differ, but are considerably different from the -0.1684 estimate that was discussed above. This confirms the main conclusion from Andrews et al. (2004) that the more movers each plant has, the smaller is the downwards bias in the correlation. Thus the estimate is a lower bound: we are not able to say whether the true correlation is zero or positive, but at least this rules out negative assortative matching.

6 AN EXAMPLE: FEMALE AND MARRIED

Here we look at the estimates on **female** (hereafter u_i), **married** (hereafter x_{it}), and their interaction most of the estimators discussed thus far. These are summarised in Table 6. Two further rows of estimates are reported, i.e. when the data are collapsed to the plant-level by taking averages over all observations in a given plant. This is as if we had plant-level data, but with observable covariates from both sides of the market. One model is Pooled OLS, the other removes plant-level heterogeneity.

Our aim here is not to attempt to contribute to the literature on the effects of marriage

¹⁹Compared with the figures given on Page 17, there are 62,668 movers, 20,313 individuals, 212 plants, and 40,719 spells.

and gender, but simply to illustrate the effects on two-well known covariates of altering the type of data observed (worker-level or plant-level) and the assumptions about unobserved worker- and plant-level heterogeneity.

The first two rows report Pooled OLS estimates; of these, the first excludes plant-level variables. Thus the second, in comparison, shows how the estimates move when data from the plant are matched to worker-level data. The estimates change very little, showing that both **female** and **married** are uncorrelated with, loosely speaking, the plant-level covariates $[\mathbf{q}_j, \mathbf{w}_{jt}]$. Using the second row, we can see that marriage is good for men, in terms of wages, by 0.0375 log-points, but is bad for women by a similar amount, namely -0.0314 log-points ($= 0.0375 - 0.0689$). Comparing these estimates with the next row, FE(i), we can see that the estimates become much smaller: 0.0058 on **married** and -0.0037 on the interaction. Comparing rows two and three suggests that marriage for men is a proxy for unobserved heterogeneity α_i and marriage for women is a proxy for $-\alpha_i$, because the returns almost disappear once α_i is controlled for. This regression, as per usual, says nothing about the effect of gender on wages (conditional on married, unmarried, or everyone), because gender is a fixed-effect.

Comparing Pooled OLS (row two) with FE(j) (row four) implies that plant-level heterogeneity is uncorrelated with either marriage or gender, because the estimates do not move. This mimics the earlier result that cross-market correlations tend to be zero. At this point we have controlled for both heterogeneity terms, but not together.

The next two rows are the two double heterogeneity methods, Spell FE and Two-step. Both methods give very similar estimates—it does not matter which technique is used—and again the estimates are small (less than one-half of one log-point). Also, they are close to the estimates that only control for worker-level heterogeneity, FE(i). An estimate on **female** can be obtained, making the usual random-effects assumption, by regressing the estimates of $\hat{\theta}_i$ on all the time-invariant individual-level variables. Not surprisingly, the estimate obtained -0.1369 is very close to the Pool OLS estimate of -0.1436 (row two). This estimate, together with 0.0056 on **married** and -0.0036 on the interaction, represent the best that can be obtained using these data.

The final two columns use the plant-level data. Here the covariates are proportion of female workers in a plant, the proportion of married workers in a plant, and the proportion of married female workers in a plant. The effect of aggregating the data is seen by comparing the second column (Pooled OLS) with the penultimate column (Plant-level, OLS), where all the estimates become much bigger in absolute values.²⁰ Wooldridge (1999) analyses exactly this situation. The estimators are unbiased if the worker-level model satisfies the Gauss-Markov assumptions and the worker-level errors are independent of plant-size. These ‘errors’ include α_i and ϕ_j . One cannot assume

²⁰It is easy to show that the effect of these estimates bear exactly the same interpretation after aggregation.

that unobservably good workers are equally likely to be employed in large and small plants; similarly one cannot assume that unobserved plant-heterogeneity is uncorrelated with plant-size. In other words, even if α_i and ϕ_j were uncorrelated with observables before aggregation, they might well be after aggregation. The second obvious source of aggregation bias occurs because α_{jt} and ϕ_j are correlated with the unobservables before aggregation and remain so afterwards. In the model being estimated:

$$y_{jt} = \mu + \mathbf{x}_{jt}\boldsymbol{\beta} + \mathbf{u}_j\boldsymbol{\eta} + \phi_j + \alpha_{jt} + \epsilon_{jt},$$

one can control for plant-level heterogeneity by using plant-level fixed effects. This makes the estimates on `female*married` much more in line with FE(j); the difference between the -0.1241 and the -0.2669 estimates on `married` is due to the fact that α_{jt} is correlated with various observables and does suggest that the true estimate is much larger than those that make the random-effects assumption. Thus one might conclude that there is some benefit to using aggregated plant-level data. If plants do not experience turnover, their α_{jt} are actually time-invariant, and get swept out, together with ϕ_j , using plant-level fixed-effects. An estimate that assumes that α_{jt} is time-invariant when it is not (for those plants that experience turnover) might be better than one that ignores α_{jt} altogether.

7 CONCLUSIONS

The main objective of this paper is to illustrate that the analysis of matched employee-employer datasets is more accessible than the investigators might imagine. We show then show how they can be implemented in Stata. We illustrate with examples using linked employer-employee data from Germany (the Linked IAB data).

There are two points worth emphasising. The first is that investigators who are interested in estimating unobserved worker heterogeneity and unobserved worker heterogeneity, and who have a ‘large’ number of plants, must use ACK’s Direct Least Squares algorithm. In this paper we explain how to make ‘large’ as small as possible—our Two Step method works well compared with the ‘correct’ FEiLSDVj method—but sometimes the regression-based techniques discussed here are not feasible.

The second point is to ask whether we actually learn we learn anything from these estimates of the worker and firm heterogeneities? It is important to emphasise that the estimates of $\hat{\psi}_j$ rely entirely on workers who change plants, as in any fixed-effects model. If one has a sample of plants, as here, there are very few movers (we have 1.9 million workers, but only 23,000 movers). The estimates on $\hat{\psi}_j$ need interpreting with caution. Moreover, we suspect that the negative correlation usually found in such studies is biased downwards, and this is caused by standard least-squares sampling

error. This issue is investigated in a companion paper.

If we do not learn anything from these estimates of the worker and firm heterogeneities, or we are not interested in them, Spell-level FE (also labelled FE(s)) is very straightforward to use.

REFERENCES

- Abowd, J., Creedy, R. & Kramarz, F. (2002), Computing person and firm effects using linked longitudinal employer-employee data, Technical Paper 2002-06, U.S. Census Bureau, April.
- Abowd, J. & Kramarz, F. (1999), The analysis of labor markets using matched employer-employee data, *in* O. Ashenfelter & D. Card, eds, 'Handbook of Labor Economics', Vol. 3B, Elsevier, Amsterdam, chapter 40, 2567–627.
- Abowd, J., Kramarz, F. & Margolis, D. (1999), 'High wage workers and high wage firms', *Econometrica* **67**, 251–333.
- Andrews, M., Schank, T. & Upward, R. (2004), High wage workers and low wage firms: negative assortative matching or statistical artefact?, Mimeo, University of Manchester, February.
- Baltagi, B. (2001), *Econometric Analysis of Panel Data*, second edn, Wiley.
- Barth, E. & Dale-Olsen, H. (2003), Assortative matching in the labour market? Stylised facts about workers and plants, Mimeo, Institute for Social Research, Oslo, February.
- Bender, S., Haas, A. & Klose, C. (2000), 'The IAB employment subsample 1975-1995: Opportunities for analysis provided by the anonymised sample', *Schmollers Jahrbuch: Zeitschrift für Wirtschafts- und Sozialwissenschaften* **120**, 649–62.
- Chamberlain, G. (1984), Panel data, *in* Z. Griliches & M. Intriligator, eds, 'Handbook of Econometrics', Vol. 2, Elsevier, Amsterdam, chapter 22, pp. 1247–318.
- Goux, D. & Maurin, E. (1999), 'Persistence of interindustry wage differentials: a reexamination using matched worker-firm panel data', *Journal of Labor Economics* **17**, 492–533.
- Gruetter, M. & Lalive, R. (2003), Job mobility and industry wage differentials: evidence from matched employer employee data, Mimeo, University of Zurich, October.
- Haltiwanger, J., Lane, J., Spletzer, J., Theeuwes, J. & Troske, K., eds (1999), *The creation and analysis of employer-employee matched data*, North-Holland.
- Hausman, J. & Taylor, W. (1981), 'Panel data and unobservable individual effects', *Econometrica* **49**, 1377–98.
- Kölling, A. (2000), 'The IAB establishment panel', *Schmollers Jahrbuch: Zeitschrift für Wirtschafts- und Sozialwissenschaften* **120**, 291–300.
- StataCorp (2003), *Stata Statistical Software: Release 8*, Stata Corporation, College Station, TX.
- Verbeek, M. (2004), *A Guide to Modern Econometrics*, second edn, Wiley.
- Wansbeek, T. & Kapteyn, A. (1989), 'Estimation of the error components models with incomplete panels', *Journal of Econometrics* **41**, 341–60.
- Wooldridge, J. (1999), *Introductory econometrics: an modern approach*, Thomson.
- Wooldridge, J. (2002), *Econometric analysis of cross section and panel data*, MIT Press.

TABLES

Table 1: The regression sample^a

Years in Sample	Number of Employers						Total	Percent
	1	1a	2	3	4	5		
1	532,875	489,896					532,875	27.6%
	1	1						
2	479,653	448,502	7,604				487,257	25.2%
	2	2	11					
3	282,599	268,095	8,102	197			290,898	15.1%
	3	3	21	111				
4	325,833	312,517	5,082	220	0		331,135	17.2%
	4	4	22	112	1111			
5	285,907	273,965	2,018	168	2	0	288,095	14.9%
	5	5	23	122	1121	11111		
Total	1,906,867	1,792,975	22,806	585	2	0	1,930,260	100.0
Percent	98.8%	92.9%	1.2%	0.0%	0.0%	0.0%	100.0%	

^aFormat of this table copied from Table 1 in Abowd et al. (1999). We report the most common employment configurations for each cell, which are described in terms of the number of consecutive years spent with each of the worker's employers (e.g. configuration 113 means that the worker spent 1 year with his first employer, then 1 year with his second employer and finally 4 years with his third employer). Column 1a refers to the subset of workers with only one employer whose employing plant had at least one other worker who had changed plants at least once in his career.

	IAB			IAB		
	all	movers	diff	all plants	turnover	diff
profit 'good' (profit2)	0.2037	0.1297	0.0739	0.2565	0.2337	0.0227
profit 'satisfactory' (profit3)	0.3772	0.4155	-0.0382	0.3550	0.3455	0.0095
profit 'just ok' (profit4)	0.2197	0.2463	-0.0266	0.2162	0.2273	-0.0112
profit 'bad' (profit5)	0.1755	0.1791	-0.0035	0.1386	0.1658	-0.0272
vtg*(1-vtgcen) ^a (vin)	0.7743	0.9064	-0.1320	2.4125	1.6574	0.7551
vin*vin (vinsq)	8.2214	5.9894	2.2320	24.1205	17.6922	6.4283
vtg*(1-vtgcen) (cvin)	15.4860	12.5803	2.9058	9.3528	12.6862	-3.3333
cvin*cvin (cvin2)	275.7671	223.0398	52.7273	1.6772	2.2738	-0.5966
1 if 1994 (year2)	0.2679	0.2438	0.0242	0.2085	0.2375	-0.0290
1 if 1995 (year3)	0.1956	0.1271	0.0685	0.2074	0.2199	-0.0124
1 if 1996 (year4)	0.1713	0.1859	-0.0147	0.2243	0.2147	0.0096
1 if 1997 (year5)	0.1505	0.1934	-0.0429	0.2163	0.1653	0.0510
No. of obs	1,930,260	23,393		4,376	1,821	

^aWhere vtg is age of the plant and vtgcen is 1 if age is censored, at 20 years.

Table 3: Conventional models^a

Pooled OLS			
	w/o [$\mathbf{q}_j, \mathbf{w}_{jt}$]	w/o [$\mathbf{u}_i, \mathbf{x}_{it}$]	
foreign	-0.0163 (0.0053)	-0.0183 (0.0038)	FE(j)
female	-0.1568 (0.0049)	-0.1436 (0.0040)	-0.0209 (0.0021)
married	0.0401 (0.0050)	0.0375 (0.0031)	-0.1241 (0.0026)
marr*fem	-0.0723 (0.0041)	-0.0689 (0.0032)	0.0419 (0.0015)
age	0.0688 (0.0033)	0.0727 (0.0029)	-0.0641 (0.0023)
age2/100	-0.1272 (0.0079)	-0.1397 (0.0068)	0.0715 (0.0027)
age3/10000	0.0775 (0.0063)	0.0893 (0.0053)	-0.1394 (0.0061)
qual2	0.1009 (0.0046)	0.0977 (0.0041)	0.0903 (0.0046)
qual3	0.1098 (0.0104)	0.0944 (0.0097)	0.0858 (0.0031)
qual4	0.1577 (0.0065)	0.1290 (0.0066)	0.0599 (0.0091)
qual5	0.2385 (0.0064)	0.2221 (0.0062)	0.1071 (0.0051)
qual6	0.2687 (0.0073)	0.2479 (0.0088)	0.1957 (0.0046)
occ2	0.0491 (0.0056)	0.0530 (0.0042)	0.2089 (0.0059)
occ3	0.2277 (0.0050)	0.2327 (0.0051)	0.0455 (0.0036)
occ4	0.0213 (0.0071)	0.0465 (0.0055)	0.2249 (0.0059)
occ5	0.1942 (0.0063)	0.1992 (0.0042)	0.0370 (0.0043)
occ6	0.2230 (0.0075)	0.2660 (0.0058)	0.1835 (0.0039)
mjob	-0.0672 (0.0090)	-0.0580 (0.0072)	0.2731 (0.0053)
single		-0.0242 (0.0082)	-0.0509 (0.0083)
region2		-0.0228 (0.0168)	-0.0081 (0.0091)
region3		-0.0424 (0.0131)	-0.0149 (0.0161)
region4		-0.1096 (0.0286)	-0.0182 (0.0095)
region5		-0.0645 (0.0171)	-0.0415 (0.0127)
region6		-0.0539 (0.0294)	-0.0445 (0.0183)
region7		-0.0984 (0.0149)	-0.0245 (0.0135)
to be continued...		-0.0515 (0.0108)	-0.0194 (0.0092)

Pooled OLS				FE(j)	
	w/o $[\mathbf{q}_j, \mathbf{w}_{jt}]$	w/o $[\mathbf{u}_i, \mathbf{x}_{it}]$		FE(i)	FE(j)
region8	-0.1059 (0.0156)	-0.0542 (0.0130)		-0.0125 (0.0100)	
region9	-0.1527 (0.0247)	-0.0884 (0.0142)		-0.0414 (0.0158)	
region10	-0.1234 (0.0437)	-0.0697 (0.0309)		-0.0275 (0.0218)	
ind2	0.1888 (0.0347)	0.1303 (0.0349)		0.0441 (0.0793)	
ind3	0.1610 (0.0307)	0.1282 (0.0326)		0.0622 (0.0776)	
ind4	0.2346 (0.0329)	0.1845 (0.0341)		-0.0016 (0.0815)	
ind5	0.0563 (0.0351)	0.0486 (0.0347)		0.0060 (0.0784)	
ind6	0.1036 (0.0326)	0.0631 (0.0341)		0.0004 (0.0815)	
ind7	0.2577 (0.0312)	0.1781 (0.0327)		0.0449 (0.0788)	
ind8	0.1258 (0.0326)	0.0517 (0.0337)		0.0039 (0.0780)	
ind9	0.0635 (0.0479)	-0.0190 (0.0419)		0.0392 (0.0803)	
ind10	0.1658 (0.0499)	0.0818 (0.0423)		0.0410 (0.0817)	
IHbar	-0.4735 (0.1448)	-0.3831 (0.0927)		-0.0605 (0.2583)	-0.0721 (0.2001)
B	-0.0343 (0.0292)	-0.0076 (0.0203)		-0.0053 (0.0093)	-0.0016 (0.0054)
B1	-0.0506 (0.0379)	-0.0245 (0.0250)		0.0004 (0.0078)	0.0085 (0.0059)
B2	-0.0481 (0.0384)	-0.0182 (0.0257)		-0.0185 (0.0102)	-0.0100 (0.0077)
inv	-0.0002 (0.0002)	-0.0002 (0.0001)		0.0001 (0.0000)	0.0001 (0.0000)
lconc	0.0005 (0.0032)	0.0035 (0.0023)		-0.0044 (0.0027)	-0.0021 (0.0024)
size1	-0.5932 (0.0351)	-0.5149 (0.0276)		-0.0818 (0.0225)	-0.0280 (0.0243)
size2	-0.4246 (0.0265)	-0.3513 (0.0213)		-0.0702 (0.0176)	-0.0030 (0.0180)
size3	-0.3038 (0.0237)	-0.2517 (0.0184)		-0.0565 (0.0163)	0.0073 (0.0156)
size4	-0.2341 (0.0201)	-0.1964 (0.0158)		-0.0455 (0.0139)	0.0143 (0.0126)
size5	-0.1992 (0.0203)	-0.1633 (0.0155)		-0.0328 (0.0115)	0.0199 (0.0111)
size6	-0.1665 (0.0173)	-0.1334 (0.0134)		-0.0228 (0.0102)	0.0117 (0.0096)
size7	-0.1217 (0.0154)	-0.0948 (0.0125)		-0.0107 (0.0096)	0.0130 (0.0087)
size8	-0.1017 (0.0157)	-0.0809 (0.0127)		-0.0039 (0.0082)	0.0072 (0.0072)
size9	-0.0681 (0.0139)	-0.0607 (0.0115)		0.0004 (0.0074)	0.0045 (0.0064)
profit2	-0.0350 (0.0153)	-0.0163 (0.0111)		-0.0015 (0.0042)	-0.0033 (0.0031)
to be continued...					

Pooled OLS				
	w/o [$\mathbf{q}_j, \mathbf{w}_{jt}$]	w/o [$\mathbf{u}_i, \mathbf{x}_{it}$]	FE(i)	FE(j)
profit3		-0.0448 (0.0168)	-0.0067 (0.0046)	-0.0082 (0.0034)
profit4		-0.0259 (0.0174)	-0.0078 (0.0051)	-0.0091 (0.0038)
profit5		-0.0440 (0.0180)	-0.0094 (0.0061)	-0.0114 (0.0045)
vin		-0.0093 (0.0052)	-0.0054 (0.0021)	0.0159 (0.0020)
vinsq		0.0003 (0.0003)	0.0001 (0.0001)	0.0002 (0.0001)
cvin		-0.0093 (0.0041)	-0.0018 (0.0020)	-0.0058 (0.0278)
cvin2		0.0004 (0.0002)	0.0000 (0.0001)	0.0007 (0.0008)
year2	0.0198 (0.0079)	0.0225 (0.0062)	-0.0047 (0.0048)	0.0012 (0.0030)
year3	0.0467 (0.0105)	0.0734 (0.0312)	0.0123 (0.0037)	0.0155 (0.0031)
year4	0.0673 (0.0105)	0.0894 (0.0318)	0.0060 (0.0023)	0.0094 (0.0018)
year5	0.0793 (0.0103)	0.1035 (0.0327)		
cons	8.3740 (0.0533)	11.5001 (0.5287)	7.5951 (0.9940)	8.5528 (0.8486)
psihat				
No. of obs	5,145,098	5,145,098	5,145,098	5,145,098
No. of workers	1,930,260	1,930,260	1,930,260	1,930,260
No. of plants	4,376	4,376	4,376	4,376
No. of spells				
'corr(ui, Xb)'	not applic	not applic	-0.6591	0.0773
σ_θ or σ_ψ	not applic	not applic	0.3529	0.2968
σ_ϵ	0.2015	0.2610	0.0680	0.1687

^a9 urbanicity dummies also included. For all regressions, we report robust standard errors adjusted for clustering on firms.

Table 4: Double heterogeneity models^a

	Spell FE	Spell FEIV	Two-step method		
			1st Step	2nd Step	thetahat psihat
foreign		-0.1091 (0.0024)			-0.1156 (0.0054)
female		-0.1360 (0.0020)			-0.1369 (0.0059)
married	0.0056 (0.0020)	0.0056 (0.0002)	0.0044 (0.0026)	0.0060 (0.0020)	
marr*fem	-0.0036 (0.0028)	-0.0036 (0.0004)	0.0003 (0.0071)	-0.0040 (0.0028)	
age	0.1035 (0.0045)	0.1035 (0.0004)	0.1142 (0.0558)	0.0895 (0.0044)	
age2/100	-0.1643 (0.0084)	-0.1643 (0.0009)	-0.2080 (0.0234)	-0.1683 (0.0085)	
age3/10000	0.1060 (0.0065)	0.1057 (0.0008)	0.1412 (0.0187)	0.1090 (0.0066)	
qual2	0.0108 (0.0053)	0.0108 (0.0008)	0.0134 (0.0073)	0.0109 (0.0049)	
qual3	-0.0615 (0.0313)	-0.0615 (0.0025)	-0.0263 (0.0188)	-0.0560 (0.0272)	
qual4	0.0157 (0.0160)	0.0157 (0.0018)	0.0150 (0.0107)	0.0176 (0.0126)	
qual5	0.0518 (0.0118)	0.0518 (0.0017)	0.0165 (0.0121)	0.0432 (0.0089)	
qual6	0.0632 (0.0158)	0.0632 (0.0020)	0.0243 (0.0134)	0.0533 (0.0121)	
occ2	0.0010 (0.0034)	0.0010 (0.0004)	0.0062 (0.0049)	0.0015 (0.0033)	
occ3	0.0359 (0.0043)	0.0359 (0.0005)	0.0428 (0.0074)	0.0360 (0.0042)	
occ4	-0.0035 (0.0031)	-0.0035 (0.0005)	-0.0103 (0.0056)	-0.0040 (0.0030)	
occ5	0.0193 (0.0045)	0.0193 (0.0006)	0.0180 (0.0086)	0.0188 (0.0045)	
occ6	0.0354 (0.0043)	0.0354 (0.0008)	0.0417 (0.0074)	0.0360 (0.0043)	
mjob	-0.0208 (0.0044)	-0.0208 (0.0015)	0.0045 (0.0231)	-0.0192 (0.0043)	
single		-0.0652 (0.0022)			-0.0328 (0.0084)
region2		-0.0210 (0.0039)			-0.0181 (0.0169)
region3		-0.0049 (0.0021)			0.0068 (0.0149)
region4		-0.0767 (0.0044)			-0.0659 (0.0404)
region5		-0.0364 (0.0053)			-0.0615 (0.0155)
region6		-0.0679 (0.0053)			-0.0354 (0.0228)
region7		-0.0629 (0.0027)			-0.0376 (0.0145)
to be continued...					

	Two-step method				
	Spell FE	Spell FEIV	1st Step	2nd Step	thetahat
profit3	-0.0072 (0.0047)	-0.0072 (0.0002)	-0.0029 (0.0096)	-0.0070 (0.0049)	psihat
profit4	-0.0082 (0.0053)	-0.0082 (0.0003)	-0.0070 (0.0065)	-0.0082 (0.0054)	
profit5	-0.0099 (0.0062)	-0.0099 (0.0003)	-0.0085 (0.0086)	-0.0094 (0.0063)	
vin			0.0184 (0.0544)	0.0156 (0.0021)	
vinsq	0.0001 (0.0001)	0.0001 (0.0000)	0.0003 (0.0003)	0.0002 (0.0001)	
cvin	-0.0160 (0.0354)	-0.0160 (0.0021)		-0.0042 (0.0022)	
cvin2	0.0005 (0.0010)	0.0005 (0.0001)	0.0005 (0.0015)	0.0006 (0.0001)	
year2	-0.0032 (0.0037)	-0.0032 (0.0002)	-0.0142 (0.0050)	-0.0036 (0.0049)	
year3	0.0140 (0.0041)	0.0140 (0.0003)	0.0227 (0.0090)	0.0149 (0.0039)	
year4	0.0073 (0.0024)	0.0073 (0.0002)	0.0091 (0.0061)	0.0076 (0.0024)	
year5					
cons	7.8941 (1.1414)	7.7591 (0.0257)		(0.0000)	8.3219 (0.0061)
psihat				0.9427 (0.0205)	-0.1262 (0.0702)
No. of obs.	5,145,098	5,145,098	72,253	4,883,331	4,883,331
No. of workers	1,930,260	1,930,260	23,393	1,816,368	1,816,368
No. of plants	4,376	4,376	1,821	1,821	1,821
No. of spells	1,953,774	1,953,774			
‘corr(ui,Xb)’	-0.5625				
σ_λ	0.3220				
σ_ϵ	0.0675		0.0692	0.0536	0.2222
					0.1122

^a9 urbanicity dummies also included. For Spell FE, 1st Step, and the auxiliary regressions for $\hat{\theta}_i$ and $\hat{\psi}_j$, we report robust standard errors adjusted for clustering on firms. Stata does not give robust standard errors for its IV GLS routine. 1st Step's standard errors have been multiplied by the scale factor 1.2237 given in Equation (8). 2nd Step's standard errors have been multiplied by the scale factor 1.2618 for the same reason.

Table 5: Models with only 212 large plants^a

	FEiLSDVj		Two-step method		
	thetahat	psihat	2nd Step	thetahat	psihat
foreign	-0.1106 (0.0076)			-0.1132 (0.0088)	
female	-0.1148 (0.0068)			-0.1001 (0.0082)	
married	0.0057 (0.0020)		0.0052 (0.0023)		
mari*fem	-0.0036 (0.0028)		-0.0041 (0.0033)		
age	0.1068 (0.0042)		0.1099 (0.0050)		
age2/100	-0.1656 (0.0082)		-0.1733 (0.0095)		
age3/10000	0.1067 (0.0064)		0.1135 (0.0073)		
qual2	0.0130 (0.0049)		0.0090 (0.0052)		
qual3	-0.0516 (0.0261)		-0.0585 (0.0295)		
qual4	0.0195 (0.0122)		0.0175 (0.0140)		
qual5	0.0459 (0.0089)		0.0427 (0.0100)		
qual6	0.0560 (0.0123)		0.0537 (0.0132)		
occ2	0.0011 (0.0033)		0.0011 (0.0034)		
occ3	0.0358 (0.0042)		0.0341 (0.0045)		
occ4	-0.0046 (0.0029)		-0.0050 (0.0031)		
occ5	0.0183 (0.0044)		0.0161 (0.0048)		
occ6	0.0348 (0.0041)		0.0330 (0.0047)		
mjob	-0.0196 (0.0042)		-0.0189 (0.0049)		
single		-0.0169 (0.0062)			-0.0120 (0.0063)
region2		0.0036 (0.0063)			0.0106 (0.0089)
region3		0.0334 (0.0161)			0.0389 (0.0160)
region4		0.0092 (0.0212)			0.0135 (0.0179)
region5		0.0005 (0.0101)			0.0021 (0.0135)
region6		-0.0046 (0.0055)			0.0102 (0.0104)
region7		-0.0016 (0.0101)			0.0065 (0.0099)
to be continued...					

FEiLSDVj			Two-step method		
	thetahat	psihat	2nd Step	thetahat	psihat
profit3	-0.0072 (0.0047)		-0.0073 (0.0056)		
profit4	-0.0081 (0.0052)		-0.0077 (0.0062)		
profit5	-0.0098 (0.0062)		-0.0088 (0.0071)		
vin	-0.0029 (0.0021)		-0.0032 (0.0023)		
vinsq	0.0001 (0.0001)		0.0002 (0.0001)		
cvin	-0.0003 (0.0020)		0.0017 (0.0024)		
cvin2	0.0000 (0.0001)		-0.0001 (0.0001)		
year2	-0.0047 (0.0048)		-0.0054 (0.0054)		
year3	0.0124 (0.0036)		0.0107 (0.0046)		
year4	0.0060 (0.0023)		0.0046 (0.0028)		
year5					
cons	7.7096 (0.0057)	-0.0066 (0.0103)	7.7588 (0.0066)	-0.0323 (0.0175)	
psihat			0.9872 (0.0348)		
No. of obs.	5,145,098	5,145,098	4,203,823	4,203,823	4,203,823
No. of workers	1,930,260	1,930,260	1,546,527	1,546,527	1,546,527
No. of plants	212	212	212	212	212
σ_ϵ	0.0675	0.3145		0.3105	0.0463

^a9 urbanicity dummies also included. For the auxiliary regressions for $\hat{\theta}_i$ and $\hat{\psi}_j$, we report robust standard errors adjusted for clustering on firms. 2nd Step's standard errors have been multiplied by the scale factor 1.2578 for the same reason.

Table 6: Female and married^a

	female	married	female*married
Pooled OLS, w/o [$\mathbf{q}_j, \mathbf{w}_{jt}$]	-0.1568 (0.0049)	0.0401 (0.0050)	-0.0723 (0.0041)
Pooled OLS	-0.1436 (0.0040)	0.0375 (0.0031)	-0.0689 (0.0032)
Within- i FE		0.0058 (0.0019)	-0.0037 (0.0027)
Within- j FE	-0.1241 (0.0026)	0.0419 (0.0015)	-0.0641 (0.0023)
Within- s FE (Spell FE)		0.0056 (0.0020)	-0.0036 (0.0028)
Two-step, 2nd step		0.0060 (0.0016)	-0.0040 (0.0022)
Two-step, auxiliary	-0.1369 (0.0059)		
Plant-level, OLS	-0.3050 (0.0245)	0.1204 (0.0200)	-0.2316 (0.0405)
Plant-level, FE	-0.2669 (0.0593)	0.0105 (0.0179)	-0.0343 (0.0432)

^aThe first seven regressions copied from Tables 3 and 4. See Section 6 for discussion of the remaining two regressions.

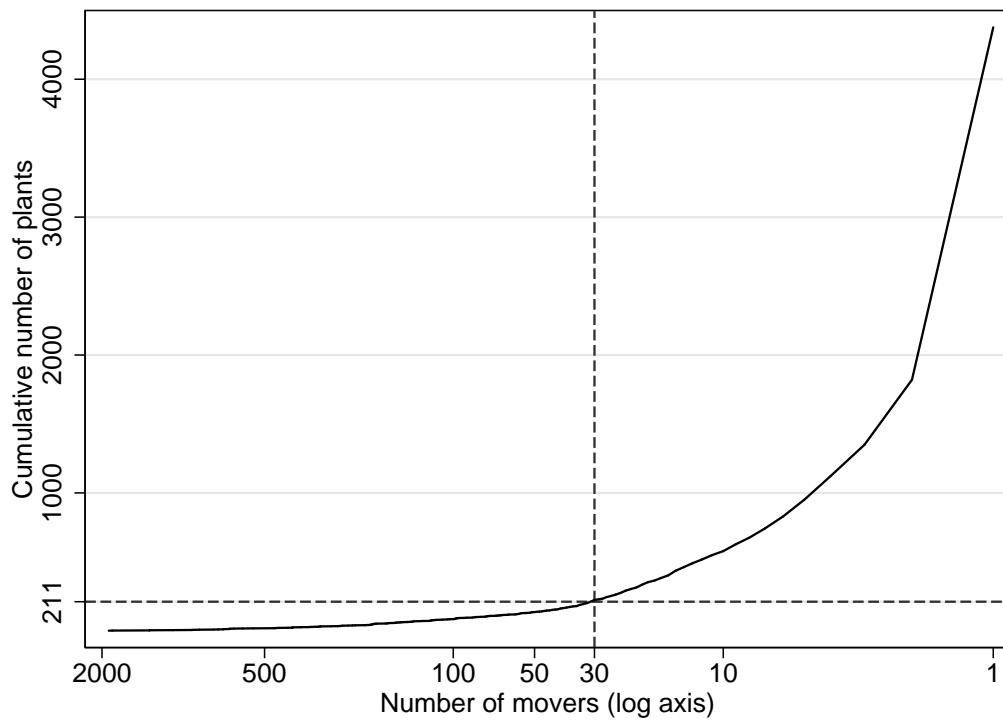
FIGURES

Figure 1: Distribution function of numbers of movers for each palnt

In der Diskussionspapierreihe sind zuletzt erschienen:

Recently published Discussion Papers:

29	Andrews, M.J., Schank, T., Upward, R.	Practical estimation methods for linked employer-employee data	09/2004
28	Brixy, U., Kohaut, S., Schnabel, C.	Do newly founded firms pay lower wages? First evidence from Germany	07/2004
27	Schank, T., Schnabel, C., Wagner, J.	Exporting firms do not pay higher wages, <i>ceteris paribus</i> . First evidence from linked employer-employee data	06/2004
26	List, J., Schnabel, C.	Bildungsstagnation bei abnehmender Erwerbsbevölkerung – Bildungspolitische Herausforderungen durch Geringsqualifizierte	05/2004
25	Andrews, M.J., Schank, T., Simmons, R.	Does Worksharing Work? Some Empirical Evidence from the IAB Panel	05/2004
24	Schank, T., Schnabel, C.	Betriebliche Determinanten des Überstundeneinsatzes	02/2004
23	Kohaut, S., Schnabel, C.	Verbreitung, Ausmaß und Determinanten der übertariflichen Entlohnung	12/2003
22	Addison, J.T., Schnabel, C., Wagner, J.	The Course of Research into the Economic Consequences of German Works Councils	11/2003
21	Addison, J.T., Schank, T., Schnabel, C., Wagner, J.	German Works Councils in the Production Process	07/2003
20	Niederalt, M.	Betriebliche Ausbildung als kollektives Phänomen	05/2003
19	Haltiwanger, J., Jarmin, R., Schank, T.	Productivity, Investment in ICT and Market Experimentation: Micro Evidence from Germany and the U.S.	03/2003

Eine aktualisierte Liste der Diskussionspapiere findet sich auf der Homepage:

<http://www.arbeitsmarkt.wiso.uni.erlangen.de/>

An updated list of discussion papers can be found at the homepage:

<http://www.arbeitsmarkt.wiso.uni.erlangen.de/>