

# An ordered probit analysis of transaction stock prices\*

Jerry A. Hausman and Andrew W. Lo

*Massachusetts Institute of Technology, Cambridge, MA 02139, USA*

A. Craig MacKinlay

*University of Pennsylvania, Philadelphia, PA 19104, USA*

Received January 1991, final version received March 1992

We estimate the **conditional distribution** of trade-to-trade price changes using **ordered probit**, a statistical model for discrete random variables. This approach recognizes that transaction price changes occur in discrete increments, typically eighths of a dollar, and occur at irregularly-spaced time intervals. Unlike existing models of discrete transaction prices, ordered probit can quantify the effects of other economic variables like volume, past price changes, and the time between trades on price changes. Using **1988 transactions data for over 100 randomly chosen U.S. stocks**, we estimate the ordered probit model via maximum likelihood and **use the parameter estimates** to measure several transaction-related quantities, such as the price impact of trades of a given size, the tendency towards price reversals from one transaction to the next, and the empirical significance of price discreteness.

## 1. Introduction

Virtually all empirical investigations of the microstructure of securities markets require a statistical model of asset prices that can capture the salient

\*We thank Arnout Eikeboom for excellent research assistance and many helpful comments, and Sarah Fisher, Ayman Hindy, and John Simpson for research assistance on earlier drafts of this paper. We have also benefited from the comments of Bill Schwert (the editor), Bruce Lehmann (the referee), Rob Bliss, Larry Harris, Joel Hasbrouck, Bob Merton, Whitney Newey, Peter Rossi, Alex Samarov, Ken Singleton, and Ronald Thisted, as well as seminar participants at Boston College, the California Institute of Technology, Carnegie Mellon University, Cornell University, the Federal Reserve Bank of Cleveland, Harvard University, Indiana University, the 1991 Johnson Symposium at the University of Wisconsin Madison, the London Business School, MIT, the NBER Summer Institute, Northwestern University, the Q Group Spring 1992 Seminar, Rice University, Stanford University, Temple University, the University of British Columbia, UC Berkeley, UCLA, the University of Chicago, the University of North Carolina, the University of Rochester, the University of Texas at Austin, Vanderbilt University, and Washington University. Research support from the Battlerymarch Fellowship, the Geewax-Terker Investments Research Fund, the MIT International Financial Services Research Center, the National Science Foundation (SES-8618769, SES-8821583), and the Q Group is gratefully acknowledged.

features of price movements from one transaction to the next. For example, because there are several theories of why bid/ask spreads exist, a stochastic model for prices is a prerequisite to empirically decomposing observed spreads into components due to order-processing costs, adverse selection, and specialist market power.<sup>1</sup> The benefits and costs of particular aspects of a market's microstructure, such as margin requirements, the degree of competition faced by dealers, the frequency that orders are cleared, and intraday volatility also depend intimately on the particular specification of price dynamics.<sup>2</sup> Even the event study, a tool that does not explicitly assume any particular theory of the market microstructure, depends heavily on price dynamics [see, for example, Barclay and Litzenberger (1988)]. In fact, it is difficult to imagine an economically relevant feature of transaction prices and the market microstructure that does *not* hinge on such price dynamics.

Since stock prices are perhaps the most closely watched economic variables to date, they have been modeled by many competing specifications, beginning with the simple random walk or Brownian motion. However, the majority of these specifications have been unable to capture at least three aspects of *transaction* prices. First, on most U.S. stock exchanges, prices are quoted in increments of eighths of a dollar, a feature not captured by stochastic processes with continuous state spaces. Of course, discreteness is less problematic for coarser-sampled data, which may be well-approximated by a continuous-state process. But discreteness is of paramount importance for intraday price movements, since such finely-sampled price changes may take on only five or six distinct values.<sup>3</sup>

The second distinguishing feature of transaction prices is their timing, which is irregular and random. Therefore, such prices may be modeled by discrete-time processes only if we are prepared to ignore the information contained in waiting times between trades.

Finally, although many studies have computed correlations between transaction price changes and other economic variables, to date none of the existing models of discrete transaction prices have been able to quantify such effects formally. Such models have focused primarily on the *unconditional* distribution of price changes, whereas what is more often of economic interest is the *conditional* distribution, conditioned on quantities such as volume, time between trades, and the *sequence* of past price changes.<sup>4</sup> For example, one of the

<sup>1</sup>See, for example, Glosten and Harris (1988), Hasbrouck (1988), Roll (1984), and Stoll (1989).

<sup>2</sup>See Cohen et al. (1986), Harris, Sofianos, and Shapiro (1990), Hasbrouck (1991a, b), Madhavan and Smidt (1991), and Stoll and Whaley (1990).

<sup>3</sup>The implications of discreteness have been considered in many studies, e.g., Cho and Frees (1988), Gottlieb and Kalay (1985), Harris (1989a, b, 1991), Petersen (1986), and Pritsker (1990).

<sup>4</sup>There is, however, a substantial literature on price/volume relations in which discreteness is ignored because of the return horizons involved (usually daily or longer). See, for example, Campbell, Grossman, and Wang (1991), Gallant, Rossi, and Tauchen (1992), and Karpoff (1987).

unresolved empirical issues in this literature is what the total costs of immediate execution are, which many take to be a measure of market liquidity. Indeed, the largest component of these costs may be the price impact of large trades. A floor broker seeking to unload 100,000 shares of stock will generally break up the sale into smaller blocks to minimize the price impact of the trades. How do we measure price impact? Such a question is a question about the conditional distribution of price changes, conditional upon a particular sequence of volume and price changes, i.e., order flow.

In this paper, we propose a specification of transaction price changes that addresses all three of these issues, and yet is still tractable enough to permit estimation via standard techniques. This specification is known as *ordered probit*, a technique used most frequently in cross-sectional studies of dependent variables that take on only a finite number of values possessing a natural ordering.<sup>5</sup> For example, the dependent variable might be the level of education, as measured by three categories: less than high school, high school, and college education. The dependent variable is discrete, and is naturally ordered since college education always follows high school [see Maddala (1983) for further details]. Heuristically, ordered probit analysis is a generalization of the linear regression model to cases where the dependent variable is discrete. As such, among the existing models of stock price discreteness [e.g., Ball (1988), Cho and Frees (1988), Gottlieb and Kalay (1985), and Harris (1991)], ordered probit is perhaps the only specification that can easily capture the impact of 'explanatory' variables on price changes while also accounting for price discreteness and irregular trade times.

Underlying the analysis is a 'virtual' regression model with an unobserved continuous dependent variable  $Z^*$  whose conditional mean is a linear function of observed 'explanatory' variables. Although  $Z^*$  is unobserved, it is related to an observable discrete random variable  $Z$ , whose realizations are determined by where  $Z^*$  lies in its domain or state space. By partitioning the state space into a finite number of distinct regions,  $Z$  may be viewed as an indicator function for  $Z^*$  over these regions. For example, a discrete random variable  $Z$  taking on the values  $\{-\frac{1}{8}, 0, \frac{1}{8}\}$  may be modeled as an indicator variable that takes on the value  $-\frac{1}{8}$  whenever  $Z^* \leq \alpha_1$ , the value 0 whenever  $\alpha_1 < Z^* \leq \alpha_2$ , and the value  $\frac{1}{8}$  whenever  $Z^* > \alpha_2$ . Ordered probit analysis consists of estimating  $\alpha_1$ ,  $\alpha_2$ , and the coefficients of the unobserved regression model that determines the conditional mean and variance of  $Z^*$ .

Since  $\alpha_1$ ,  $\alpha_2$ , and  $Z^*$  may depend on a vector of 'regressors'  $X$ , ordered probit analysis is considerably more general than its simple structure suggests. In fact, it is well-known that ordered probit can fit any arbitrary multinomial

<sup>5</sup>The ordered probit model was developed by Aitchison and Silvey (1957) and Ashford (1959), and generalized to nonnormal disturbances by Gurland, Lee, and Dahm (1960). For more recent extensions, see Maddala (1983), McCullagh (1980), and Thisted (1991).

distribution. However, because of the underlying linear regression framework, ordered probit can also capture the price effects of many economic variables in a way that models of the unconditional distribution of price changes cannot.

To motivate our methodology and to focus it on specific market microstructure applications, we consider three questions concerning the behavior of transaction prices. First, how does the particular sequence of trades affect the conditional distribution of price changes, and how do these effects differ across stocks? For example, does a sequence of three consecutive buyer-initiated trades ['buys'] generate price pressure, so that the next price change is more likely to be positive than if the sequence were three consecutive seller-initiated trades ['sells'], and how does this pressure change from stock to stock? Second, does trade size affect price changes as some theories suggest, and if so, what is the price impact per unit volume of trade from one transaction to the next? Third, does price discreteness matter? In particular, can the conditional distribution of price changes be modeled as a simple linear regression of price changes on explanatory variables without accounting for discreteness at all? Within the context of the ordered probit framework, we shall obtain sharp answers to each of these questions.

In section 2, we review the ordered probit model and describe its estimation via maximum likelihood. We describe the data in section 3 by presenting detailed summary statistics for an initial sample of six stocks. In section 4, we discuss the empirical specification of the ordered probit model and the selection of conditioning or 'explanatory' variables. The maximum likelihood estimates for our initial sample are reported in section 5, along with some diagnostic specification tests. In section 6, we use these maximum likelihood estimates in three specific applications: (1) testing for order-flow dependence, (2) measuring price impact, and (3) comparing ordered probit to simple linear regression. And as a check on the robustness of our findings, in section 7 we present less detailed results for a larger and randomly chosen sample of 100 stocks. We conclude in section 8.

## **2. The ordered probit model**

Consider a sequence of transaction prices  $P(t_0), P(t_1), P(t_2), \dots, P(t_n)$  observed at times  $t_0, t_1, t_2, \dots, t_n$ , and denote by  $Z_1, Z_2, \dots, Z_n$  the corresponding price changes, where  $Z_k \equiv P(t_k) - P(t_{k-1})$  is assumed to be an integer multiple of some divisor called a 'tick' [such as an eighth of a dollar]. Let  $Z_k^*$  denote an unobservable continuous random variable such that

$$Z_k^* = X_k' \beta + \varepsilon_k, \quad E[\varepsilon_k | X_k] = 0, \quad \varepsilon_k \text{ i.n.i.d. } N(0, \sigma_k^2), \quad (2.1)$$

where ‘i.n.i.d.’ indicates that the  $\varepsilon_k$ ’s are independently but *not* identically distributed, and  $X_k$  is a  $q \times 1$  vector of predetermined variables that governs the conditional mean of  $Z_k^*$ . Note that subscripts are used to denote ‘transaction’ time, whereas time arguments  $t_k$  denote calendar or ‘clock’ time, a convention we shall follow throughout the paper.

The essence of the ordered probit model is the assumption that observed price changes  $Z_k$  are related to the continuous variable  $Z_k^*$  in the following manner:

$$Z_k = \begin{cases} s_1 & \text{if } Z_k^* \in A_1, \\ s_2 & \text{if } Z_k^* \in A_2, \\ \vdots & \vdots \\ s_m & \text{if } Z_k^* \in A_m, \end{cases} \quad (2.2)$$

where the sets  $A_j$  form a *partition* of the state space  $\mathcal{S}^*$  of  $Z_k^*$ , i.e.,  $\mathcal{S}^* = \bigcup_{j=1}^m A_j$  and  $A_i \cap A_j = \emptyset$  for  $i \neq j$ , and the  $s_j$ ’s are the discrete values that comprise the state space  $\mathcal{S}$  of  $Z_k$ .

The motivation for the ordered probit specification is to uncover the mapping between  $\mathcal{S}^*$  and  $\mathcal{S}$  and relate it to a set of economic variables or ‘regressors’. In our current application, the  $s_j$ ’s are  $0, -\frac{1}{8}, +\frac{1}{8}, -\frac{2}{8}, +\frac{2}{8}$ , and so on, and for simplicity we define the state-space partition of  $\mathcal{S}^*$  to be intervals:

$$A_1 \equiv (-\infty, \alpha_1], \quad (2.3)$$

$$A_2 \equiv (\alpha_1, \alpha_2], \quad (2.4)$$

$$\vdots$$

$$A_i \equiv (\alpha_{i-1}, \alpha_i], \quad (2.5)$$

$$\vdots$$

$$A_m \equiv (\alpha_{m-1}, \infty). \quad (2.6)$$

Although the observed price change can be any number of ticks, positive or negative, we assume that  $m$  in (2.2) is finite to keep the number of unknown parameters finite. This poses no problems, since we may always let some states in  $\mathcal{S}$  represent a multiple [and possibly countably infinite] number of values for the observed price change. For example, in our empirical application we define  $s_1$  to be a price change of  $-4$  ticks *or less*,  $s_9$  to be a price change of  $+4$  ticks *or more*, and  $s_2$  to  $s_8$  to be price changes of  $-3$  ticks to  $+3$  ticks, respectively. This parsimony is obtained at the cost of *losing price resolution* — under this specification the ordered probit model does not distinguish between price changes of  $+4$  and price changes greater than  $+4$  [since the  $+4$ -tick outcome and the greater than  $+4$ -tick outcome have been grouped together into

a common event], and similarly for price changes of  $-4$  ticks versus price changes less than  $-4$ . Of course, in principle the resolution may be made arbitrarily finer by simply introducing more states, i.e., by increasing  $m$ . Moreover, as long as (2.1) is correctly specified, then increasing price resolution will not affect the estimated  $\beta$ 's asymptotically [although finite sample properties may differ]. However, in practice the data will impose a limit on the fineness of price resolution simply because there will be no observations in the extreme states when  $m$  is too large, in which case a subset of parameters is not identified and cannot be estimated.

Observe that the  $\varepsilon_k$ 's in (2.1) are assumed to be conditionally independently but *not* identically distributed, conditioned on the  $X_k$ 's and other economic variables  $W_k$  influencing the conditional variance  $\sigma_k^2$ .<sup>6</sup> This allows for clock-time effects, as in the case of an arithmetic Brownian motion where the variance  $\sigma_k^2$  of price changes is linear in the time between trades. We also allow for more general forms of conditional heteroskedasticity by letting  $\sigma_k^2$  depend linearly on other economic variables  $W_k$ , which differs from Engle's (1982) ARCH process only in its application to a discrete dependent variable model requiring an additional identification assumption that we shall discuss below in section 4.

The dependence structure of the observed process  $Z_k$  is clearly induced by that of  $Z_k^*$  and the definitions of the  $A_j$ 's, since

$$P(Z_k = s_j | Z_{k-1} = s_i) = P(Z_k^* \in A_j | Z_{k-1}^* \in A_i). \quad (2.7)$$

As a consequence, if the variables  $X_k$  and  $W_k$  are temporally independent, the observed process  $Z_k$  is also temporally independent. Of course, these are fairly restrictive assumptions and are certainly not necessary for any of the statistical inferences that follow. We require only that the  $\varepsilon_k$ 's be *conditionally* independent, so that all serial dependence is captured by the  $X_k$ 's and the  $W_k$ 's. Consequently, the independence of the  $\varepsilon_k$ 's does not imply that the  $Z_k^*$ 's are independently distributed because we have placed no restrictions on the temporal dependence of the  $X_k$ 's or  $W_k$ 's.

The conditional distribution of observed price changes  $Z_k$ , conditioned on  $X_k$  and  $W_k$ , is determined by the partition boundaries and the particular distribution of  $\varepsilon_k$ . For Gaussian  $\varepsilon_k$ 's, the conditional distribution is

$$\begin{aligned} & P(Z_k = s_i | X_k, W_k) \\ &= P(X_k' \beta + \varepsilon_k \in A_i | X_k, W_k) \end{aligned} \quad (2.8)$$

<sup>6</sup>Unless explicitly stated otherwise, all the probabilities we deal with in this study are conditional probabilities, and all statements concerning these probabilities are conditional statements, conditioned on these variables.

$$= \begin{cases} P(X_k' \beta + \varepsilon_k \leq \alpha_1 | X_k, W_k) & \text{if } i = 1, \\ P(\alpha_{i-1} < X_k' \beta + \varepsilon_k \leq \alpha_i | X_k, W_k) & \text{if } 1 < i < m, \\ P(\alpha_{m-1} < X_k' \beta + \varepsilon_k | X_k, W_k) & \text{if } i = m, \end{cases} \quad (2.9)$$

$$= \begin{cases} \Phi\left(\frac{\alpha_1 - X_k' \beta}{\sigma_k}\right) & \text{if } i = 1, \\ \Phi\left(\frac{\alpha_i - X_k' \beta}{\sigma_k}\right) - \Phi\left(\frac{\alpha_{i-1} - X_k' \beta}{\sigma_k}\right) & \text{if } 1 < i < m, \\ 1 - \Phi\left(\frac{\alpha_{m-1} - X_k' \beta}{\sigma_k}\right) & \text{if } i = m, \end{cases} \quad (2.10)$$

where  $\Phi(\cdot)$  is the standard normal cumulative distribution function.

To develop some intuition for the ordered probit model, observe that the probability of any particular observed price change is determined by where the conditional mean lies relative to the partition boundaries. Therefore, for a given conditional mean  $X_k' \beta$ , shifting the boundaries will alter the probabilities of observing each state [see fig. 1]. In fact, by shifting the boundaries

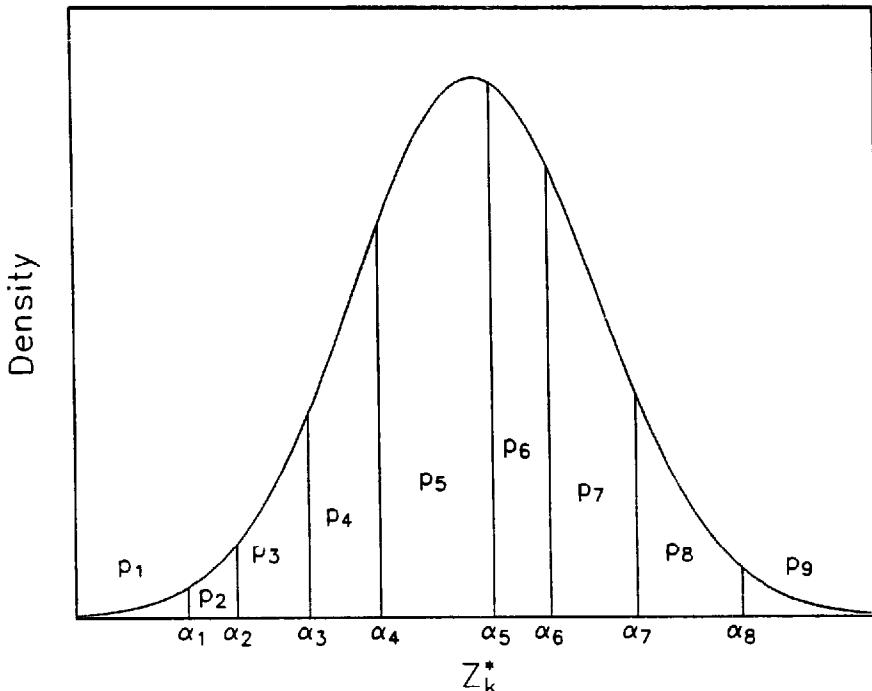


Fig. 1. Illustration of ordered probit probabilities  $p_i$  of observing a price change of  $s_i$  ticks, which are determined by where the unobservable 'virtual' price change  $Z_k^*$  falls. In particular, if  $Z_k^*$  falls in the interval  $(\alpha_{i-1}, \alpha_i]$ , then the ordered probit model implies that the observed price change  $Z_k$  is  $s_i$  ticks. More formally,  $p_i \equiv \text{Prob}(Z_k = s_i | X_k, W_k) = \text{Prob}(\alpha_{i-1} < Z_k^* \leq \alpha_i | X_k, W_k)$ ,  $i = 1, \dots, 9$ , where, for notational simplicity, we define  $\alpha_0 \equiv -\infty$  and  $\alpha_9 \equiv +\infty$ . The ordered probit model captures the effect of economic variables  $X_k, W_k$  on the virtual price change and places enough structure on the probabilities  $p_i$  to permit their estimation by maximum likelihood.

appropriately, ordered probit can fit any arbitrary multinomial distribution. This implies that the assumption of normality underlying ordered probit plays no special role in determining the probabilities of states; a logistic distribution, for example, could have served equally well. However, since it is considerably more difficult to capture conditional heteroskedasticity in the ordered logit model, we have chosen the Gaussian specification.

Given the partition boundaries, a higher conditional mean  $X_k'\beta$  implies a higher probability of observing a more extreme positive state. Of course, the labeling of states is arbitrary, but the *ordered* probit model makes use of the natural ordering of the states. The regressors allow us to separate the effects of various economic factors that influence the likelihood of one state versus another. For example, suppose that a large positive value of  $X_1$  usually implies a large negative observed price change and vice versa. Then the ordered probit coefficient  $\beta_1$  will be negative in sign and large in magnitude [relative to  $\sigma_k$  of course].

By allowing the data to determine the partition boundaries  $\alpha$ , the coefficients  $\beta$  of the conditional mean, and the conditional variance  $\sigma_k^2$ , the ordered probit model captures the empirical relation between the unobservable continuous state space  $\mathcal{S}^*$  and the observed discrete state space  $\mathcal{S}$  as a function of the economic variables  $X_k$  and  $W_k$ .

### *2.1. Other models of discreteness*

From these observations, it is apparent that the rounding/eighths-barriers models of discreteness in Ball (1988), Cho and Frees (1988), Gottlieb and Kalay (1985), and Harris (1991) may be reparametrized as ordered probit models. Consider first the case of a ‘true’ price process that is an arithmetic Brownian motion, with trades occurring only when this continuous-state process crosses an eighths threshold [see Cho and Frees (1988)]. Observed trades from such a process may be generated by an ordered probit model in which the partition boundaries are fixed at multiples of eighths and the single regressor is the time interval [or first-passage time] between crossings, appearing in both the conditional mean and variance of  $Z_k^*$ .

To obtain the rounding models of Ball (1988), Gottlieb and Kalay (1985), and Harris (1991), which do not make use of waiting times between trades, define the partition boundaries as the midpoint between eighths, e.g., the observed price change is  $\frac{3}{8}$  if the virtual price process lies in the interval  $[\frac{5}{16}, \frac{7}{16}]$ , and omit the waiting time as a regressor in both the conditional mean and variance [see the discussion in section 6.3 below].

The generality of the ordered probit model comes from the fact that the rounding and eighths-barriers models of discreteness are both special cases of ordered probit under appropriate restrictions on the partition boundaries. In fact, since the boundaries may be parametrized to be time- and state-dependent,

ordered probit can allow for considerably more general kinds of rounding and eighths barriers. In addition to fitting any arbitrary multinomial distribution, ordered probit may also accommodate finite-state Markov chains and compound Poisson processes.

Of course, other models of discreteness are not necessarily obsolete, since in several cases the parameters of interest may not be simple functions of the ordered probit parameters. For example, a tedious calculation will show that although Harris's (1991) rounding model may be represented as an ordered probit model, the bid/ask spread parameter  $c$  is not easily recoverable from the ordered probit parameters. In such cases, other equivalent specifications may allow more direct estimation of the parameters of interest.

## 2.2. The likelihood function

Let  $Y_{ik}$  be an indicator variable which takes on the value one if the realization of the  $k$ th observation  $Z_k$  is the  $i$ th state  $s_i$ , and zero otherwise. Then the log-likelihood function  $\mathcal{L}$  for the vector of price changes  $Z = [Z_1 \ Z_2 \ \dots \ Z_n]'$ , conditional on the explanatory variables  $X = [X_1 \ X_2 \ \dots \ X_n]'$ , is given by

$$\begin{aligned} \mathcal{L}(Z | X) = & \sum_{k=1}^n \left\{ Y_{1k} \cdot \log \Phi \left( \frac{\alpha_1 - X'_k \beta}{\sigma_k} \right) \right. \\ & + \sum_{i=2}^{m-1} Y_{ik} \cdot \log \left[ \Phi \left( \frac{\alpha_i - X'_k \beta}{\sigma_k} \right) - \Phi \left( \frac{\alpha_{i-1} - X'_k \beta}{\sigma_k} \right) \right] \\ & \left. + Y_{mk} \cdot \log \left[ 1 - \Phi \left( \frac{\alpha_{m-1} - X'_k \beta}{\sigma_k} \right) \right] \right\}. \end{aligned} \quad (2.11)$$

Recall that  $\sigma_k^2$  is a conditional variance, conditioned upon  $X_k$ . This allows for conditional heteroskedasticity in the  $Z_k^*$ 's, as in the rounding model of Cho and Frees (1988) where the  $Z_k^*$ 's are increments of arithmetic Brownian motion with variance proportional to  $t_k - t_{k-1}$ . In fact, arithmetic Brownian motion may be accommodated explicitly by the specification

$$X'_k \beta = \mu \Delta t_k, \quad (2.12)$$

$$\sigma_k^2 = \gamma^2 \Delta t_k. \quad (2.13)$$

More generally, we may also let  $\sigma_k^2$  depend on other economic variables  $W_k$ , so that

$$\sigma_k^2 = \gamma_0^2 + \sum_{i=1}^{K_a} \gamma_i^2 W_{ik}. \quad (2.14)$$

There are, however, some constraints that must be placed on these parameters to achieve identification since, for example, doubling the  $\alpha$ 's, the  $\beta$ 's, and  $\sigma_k$  leaves the likelihood unchanged. We shall return to this issue in section 4.

### **3. The data**

The Institute for the Study of Security Markets [ISSM] transaction database consists of time-stamped trades [to the nearest second], trade size, and bid/ask quotes from the New York and American Stock Exchanges and the consolidated regional exchanges from January 4 to December 30 of 1988. Because of the sheer size of the ISSM database, most empirical studies of the market microstructure have concentrated on more manageable subsets of the database, and we also follow this practice. But because there is so much data, the 'pretest' or 'data-snooping' biases associated with any nonrandom selection procedure used to obtain the smaller subsets are likely to be substantial. As a simple example of such a bias, suppose we choose our stocks by the seemingly innocuous requirement that they have a minimum of 100,000 trades in 1988. This rule will impart a substantial downward bias on our measures of price impact because stocks with over 100,000 trades per year are generally more liquid and, almost by definition, have smaller price impact. Therefore, how we choose our subsample of stocks may have important consequences for how our results are to be interpreted, so we shall describe our procedure in some detail here.

We first begin with an initial 'test' sample containing five stocks that did not engage in any stock splits or stock dividends greater than 3:2 during 1988: Alcoa, Allied Signal, Boeing, DuPont, and General Motors. We restrict splits because the effects of price discreteness to be captured by our model are likely to change in important ways with dramatic shifts in the price level; by eliminating large splits we reduce the problem of large changes in the price level without screening on prices directly. [Of course, if we were interested in explaining stock splits, this procedure would obviously impart important biases in the empirical results.] We also chose these five stocks because they are relatively large and visible companies, each with a large number of trades, and therefore likely to yield accurate parameter estimates. We then performed the standard 'specification searches' on these five stocks, adding, deleting, and transforming regressors to obtain a 'reasonable' fit. By 'reasonable' we mean primarily the convergence of the maximum likelihood estimation procedure, but it must also include Leamer's (1978) kind of informal or *ad hoc* inferences that all empiricists engage in.

Once we obtain a specification that is 'reasonable', we estimate it *without further revision* for our primary sample of six new stocks, chosen to yield a representative sample with respect to industries, market value, price levels, and

sample sizes. They are International Business Machines Corporation (IBM), Quantum Chemical Corporation (CUE), Foster Wheeler Corporation (FWC), Handy and Harman Company (HNH), Navistar International Corporation (NAV), and American Telephone and Telegraph Incorporated (T). [Our original primary sample consists of eleven stocks but we omitted the results for five of them to conserve space. See Hausman, Lo, and MacKinlay (1991) for the full set of results.] By using the specification derived from the test sample on stocks in this fresh sample, we seek to lessen the impact of any data-snooping biases generated by our specification searches. If, for example, our parameter estimates and subsequent inferences change dramatically in the new sample [in fact, they do not] this might be a sign that our test-sample findings were driven primarily by selection biases.

As a final check on the robustness of our specification, we estimate it for a larger sample of 100 stocks chosen randomly, and these companies are listed in table 5. From this larger sample, it is apparent that our smaller six-stock sample does suffer from at least one selection bias: it is comprised of relatively well-known companies. In contrast, relatively few companies in table 5 are as familiar. Despite this bias, virtually all of our empirical findings are confirmed by the larger sample. To conserve space and to focus attention on our findings, we report the complete set of summary statistics and estimation results only for the smaller sample of six stocks, and present broader and less detailed findings for the extended sample afterwards.

Of course, as long as there is cross-sectional dependence among the two samples it is impossible to eliminate such biases completely. Moreover, samples drawn from a different time period are not necessarily free from selection bias as some have suggested, due to the presence of *temporal* dependence. Unfortunately, nonexperimental inference is always subject to selection biases of one kind or another since specification searches are an unavoidable aspect of genuine progress in empirical research [see, for example, Lo and MacKinlay (1990b)]. Even Bayesian inference, which is not as sensitive to the kinds of selection biases discussed in Leamer (1978), can be distorted in subtle ways by specification searches. Therefore, beyond our test-sample procedure, we can only alert readers to the possibility of such biases and allow them to adjust their own inferences accordingly.

### 3.1. Sample statistics

We take as our basic time series the *intraday* price changes from trade to trade, and discard all overnight price changes. That the statistical properties of overnight price changes differ considerably from those of intraday price changes has been convincingly documented by several authors, most recently by Amihud and Mendelson (1987), Stoll and Whaley (1990), and Wood, McInish, and Ord (1985). Since the three market microstructure applications we are focusing on

involve intraday price behavior, and overnight price changes are different enough to warrant a separate specification, we use only intraday price changes. The first and last transaction prices of each day are also discarded, since they differ systematically from other prices due to institutional features [see Amihud and Mendelson (1987) for further details].

Several other screens were imposed to eliminate 'problem' trades and quotes, yielding sample sizes ranging from 3,174 trades for HNH to 206,794 trades for IBM. Specifically: (1) all trades flagged with the following ISSM condition codes were eliminated: A, C, D, O, R, and Z [see the ISSM documentation for further details concerning trade condition codes]; (2) transactions exceeding 3,276,000 shares [termed 'big trades' by ISSM] were also eliminated; (3) because we use three lags of price changes and three lags of five-minute returns on the S&P500 index futures prices as explanatory variables, we do not use the first three price changes or price changes during the first 15 minutes of each day [whichever occurs later] as observations of the dependent variable; and (4) since S&P500 futures data were not available on November 10, 11, and the first two trading hours of May 3, trades during these times were also omitted.

For some stocks, a small number of transactions occurred at prices denominated in  $\frac{1}{16}$ 's,  $\frac{1}{32}$ 's, or  $\frac{1}{64}$ 's of a dollar [non-NYSE trades]. In these cases, we rounded the price randomly [up or down] to the nearest  $\frac{1}{8}$ , and if necessary, also rounded the bid/ask quotes in the same direction.

Quotes implying bid/ask spreads greater than 40 ticks or flagged with the following ISSM condition codes were also eliminated: C, D, F, G, I, L, N, P, S, V, X, and Z [essentially all 'BBO-ineligible' quotes; see the ISSM documentation for further details concerning the definitions of the particular trade and quote condition codes, and Eikeboom (1992) for a thorough study of the relative frequencies of these condition codes for a small subset of the ISSM database].

Since we also use bid and ask prices in our analysis, some discussion of how we matched quotes to prices is necessary. Bid/ask quotes are reported on the ISSM tape only when they are revised, hence it is natural to match each transaction price to the most recently reported quote *prior* to the transaction. However, Bronfman (1991), Lee and Ready (1991), and others have shown that prices of trades that precipitate quote revisions are sometimes reported with a lag, so that the order of quote revision and transaction price is reversed in official records such as the ISSM tapes. To address this issue, we match transaction prices to quotes that are set *at least five seconds prior* to the transaction; the evidence in Lee and Ready (1991) suggests that this will account for most of the missequencing.

To provide some intuition for this enormous dataset, we report a few summary statistics in table 1. Our sample contains considerable price dispersion, with the low stock price ranging from \$3.125 for NAV to \$104.250 for IBM, and the high ranging from \$7.875 for NAV to \$129.500 for IBM. At \$219 million,

Table 1

Summary statistics for transaction prices and corresponding ordered probit explanatory variables of International Business Machines Corporation (IBM – 206,794 trades), Quantum Chemical Corporation (CUE – 26,927 trades), Foster Wheeler Corporation (FWC – 18,199 trades), Handy and Harman Company (HNH – 3,174 trades), Navistar International Corporation (NAV – 96,127 trades), and American Telephone and Telegraph Company (T – 180,726 trades), for the period from 4 January 1988 to 30 December 1988.

Statistic	IBM	CUE	FWC	HNH	NAV	T
Low price	104.250	65.500	11.500	14.250	3.125	24.125
High price	129.500	108.250	17.250	18.500	7.875	30.375
Market value (\$billions) <sup>a</sup>	69.815	2.167	0.479	0.219	0.998	28.990
% trades at prices:						
> Midquote	43.81	43.19	37.13	22.53	40.80	32.3%
= Midquote	12.66	18.67	23.58	26.28	18.11	25.92
< Midquote	43.53	38.14	39.29	51.20	41.09	41.71
Price change, $Z_k$						
Mean	-0.0010	0.0016	-0.0017	-0.0028	-0.0002	0.0001
Std. dev.	0.7530	1.2353	0.6390	0.7492	0.6445	0.6540
Time between trades, $\Delta t_k$						
Mean	27.21	203.52	296.54	1129.37	58.36	31.00
Std. dev.	34.13	282.16	416.49	1497.44	76.53	34.39
Bid/ask spread, $AB_k$						
Mean	1.9470	3.2909	2.0830	2.4707	1.4616	1.6564
Std. dev.	1.4625	1.6203	1.1682	0.8994	0.6713	0.7936
S&P500 futures return, $SP500_k^b$						
Mean	-0.0000	-0.0004	-0.0017	-0.0064	0.0001	-0.0001
Std. dev.	0.0716	0.1387	0.1475	0.1963	0.1038	0.0765
Buy/sell indicator, $IBS_k^c$						
Mean	0.0028	0.0505	-0.0216	-0.2867	-0.0028	-0.0933
Std. dev.	0.9346	0.9005	0.8739	0.8095	0.9049	0.8556
Signed transformed volume <sup>d</sup>						
Mean	0.1059	0.3574	-0.0523	-1.9543	0.0332	-0.4256
Std. dev.	6.1474	5.6643	6.2798	6.0890	6.9705	7.5846
Median trading volume (\$)	57,375	40,900	6,150	5,363	3,000	7,950

<sup>a</sup>Computed at the beginning of the sample period.

<sup>b</sup>Five-minute continuously-compounded returns of the S&P500 index futures price, for the contract maturing in the closest month beyond the month in which transaction  $k$  occurred, where the return corresponding to the  $k$ th transaction is computed with the futures price recorded one minute before the nearest round minute *prior* to  $t_k$  and the price recorded five minutes before this.

<sup>c</sup>Takes the value 1 if the  $k$ th transaction price is greater than the average of the quoted bid and ask prices at time  $t_k$ , the value -1 if the  $k$ th transaction price is less than the average of the quoted bid and ask prices at time  $t_k$ , and 0 otherwise.

<sup>d</sup>Box-Cox transformation of dollar volume multiplied by the buy/sell indicator, where the Box-Cox parameter  $\lambda$  is estimated jointly with the other ordered probit parameters via maximum likelihood. The Box-Cox parameter  $\lambda$  determines the degree of curvature that the transformation  $T_\lambda(\cdot)$  exhibits in transforming dollar volume  $V_k$  before inclusion as an explanatory variable in the ordered probit specification. If  $\lambda = 1$ , the transformation  $T_\lambda(\cdot)$  is linear, hence dollar volume enters the ordered probit model linearly. If  $\lambda = 0$ , the transformation is equivalent to  $\log(\cdot)$ , hence the natural logarithm of dollar volume enters the ordered probit model. When  $\lambda$  is between 0 and 1, the curvature of  $T_\lambda(\cdot)$  is between logarithmic and linear.

HNH has the smallest market capitalization in our sample, and IBM has the largest with a market value of \$69.8 billion.

For our empirical analysis we also require some indicator of whether a transaction was buyer-initiated or seller-initiated. Obviously, this is a difficult task because for every trade there is always a buyer and a seller. What we are attempting to measure is which of the two parties is more anxious to consummate the trade and is therefore willing to pay for it in the form of the bid/ask spread. Perhaps the most obvious indicator is whether the transaction occurs at the ask price or at the bid price; if it is the former then the transaction is most likely a 'buy', and if it is the latter then the transaction is most likely a 'sell'. Unfortunately, a large number of transactions occur at prices strictly *within* the bid/ask spread, so that this method for signing trades will leave the majority of trades indeterminate.

Following Blume, MacKinlay, and Terker (1989) and many others, we classify a transaction as a buy if the transaction price is higher than the mean of the prevailing bid/ask quote [the most recent quote that is set at least five seconds prior to the trade], and classify it as a sell if the price is lower. Should the price equal the mean of the prevailing bid/ask quote, we classify the trade as an 'indeterminate' trade. This method yields far fewer indeterminate trades than classifying according to transactions at the bid or at the ask.

Unfortunately, little is known about the relative merits of this method of classification versus others such as the 'tick test' [which classifies a transaction as a buy, a sell, or indeterminate if its price is greater than, less than, or equal to the previous transaction's price, respectively], simply because it is virtually impossible to obtain the data necessary to evaluate these alternatives. The only study we have seen is by Robinson (1988, ch. 4.4.1, table 19), in which he compared the tick test rule to the bid/ask mean rule for a sample of 196 block trades initiated by two major Canadian life insurance companies, and concluded that the bid/ask mean rule was considerably more accurate.

From table 1 we see that 13–26% of each stock's transactions are indeterminate, and the remaining trades fall almost equally into the two remaining categories. The one exception is the smallest stock, HNH, which has more than twice as many sells as buys.

The means and standard deviations of other variables to be used in our ordered probit analysis are also given in table 1. The precise definitions of these variables will be given below in section 4, but briefly,  $Z_k$  is the price change between transactions  $k - 1$  and  $k$ ,  $\Delta t_k$  is the time elapsed between these trades,  $AB_k$  is the bid/ask spread prevailing at transaction  $k$ ,  $SP500_k$  is the return on the S&P500 index futures price over the five-minute period immediately preceding transaction  $k$ ,  $IBS_k$  is the buy/sell indicator described above [1 for a buy, -1 for a sell, and 0 for an indeterminate trade], and  $T_\lambda(V_k)$  is a transformation of the dollar volume of transaction  $k$ , transformed according to the Box and Cox (1964) specification with parameter  $\lambda_i$  which is estimated for each stock  $i$  by maximum likelihood along with the other ordered probit parameters.

From table 1 we see that for the larger stocks, trades tend to occur almost every minute on average. Of course, the smaller stocks trade less frequently, with HNH trading only once every 18 minutes on average. The median dollar volume per trade also varies considerably, ranging from \$3,000 for relatively low-priced NAV to \$57,375 for higher-priced IBM.

Finally, fig. 2 contains histograms for the price change, time-between-trade, and dollar volume variables for the six stocks. The histograms of price changes are constructed so that the most extreme cells also include observations *beyond* them, i.e., the level of the histogram for the - 4 tick cell reflects all price changes of - 4 ticks or less, and similarly for the + 4 ticks cell. Surprisingly, these price histograms are remarkably symmetric across all stocks. Also, virtually all the mass in each histogram is concentrated in five or seven cells – there are few absolute price changes of four ticks or more, which underscores the importance of discreteness in transaction prices.

For both the time-between-trade and dollar volume variables, the *largest* cell, i.e., 1,500 seconds or \$200,000, also includes all trades beyond it. As expected, the histograms for these quantities vary greatly according to market value and price level. For the larger stocks, the time between trades is relatively short, hence most of the mass in those histograms is in the lower-valued cells. But the histograms of smaller, less liquid stocks like HNH have spikes in the largest-valued cell. Histograms for dollar volume are sometimes bimodal, as in the case of IBM, reflecting both round-lot trading at 100 shares [\$10,000 on average for IBM's stock price during 1988] and some very large trades, presumably by institutional investors.

#### 4. The empirical specification

To estimate the parameters of the ordered probit model via maximum likelihood, we must first specify (i) the number of states  $m$ , (ii) the explanatory variables  $X_k$ , and (iii) the parametrization of the variance  $\sigma_k^2$ .

**In choosing  $m$ , we must balance price resolution against the practical constraint that too large an  $m$  will yield no observations in the extreme states  $s_1$  and  $s_m$ .** For example, if we set  $m$  to 101 and define the states  $s_1$  and  $s_{101}$  symmetrically to be price changes of - 50 ticks and + 50 ticks respectively, we would find no  $Z_k$ 's among our six stocks falling into these two states. Using the histograms in fig. 2 as a guide, we set  $m = 9$  for the larger stocks, implying extreme states of - 4 ticks or less and + 4 ticks or more. For the two smaller stocks, FWC and HNH, we set  $m = 5$ , implying extreme states of - 2 ticks or less and + 2 ticks or more. Although the definition of states need not be symmetric [state  $s_1$  can be - 6 ticks or less, implying that state  $s_9$  is + 2 ticks or more], the symmetry of the histogram of price changes in fig. 2 suggests a symmetric definition of the  $s_j$ 's.

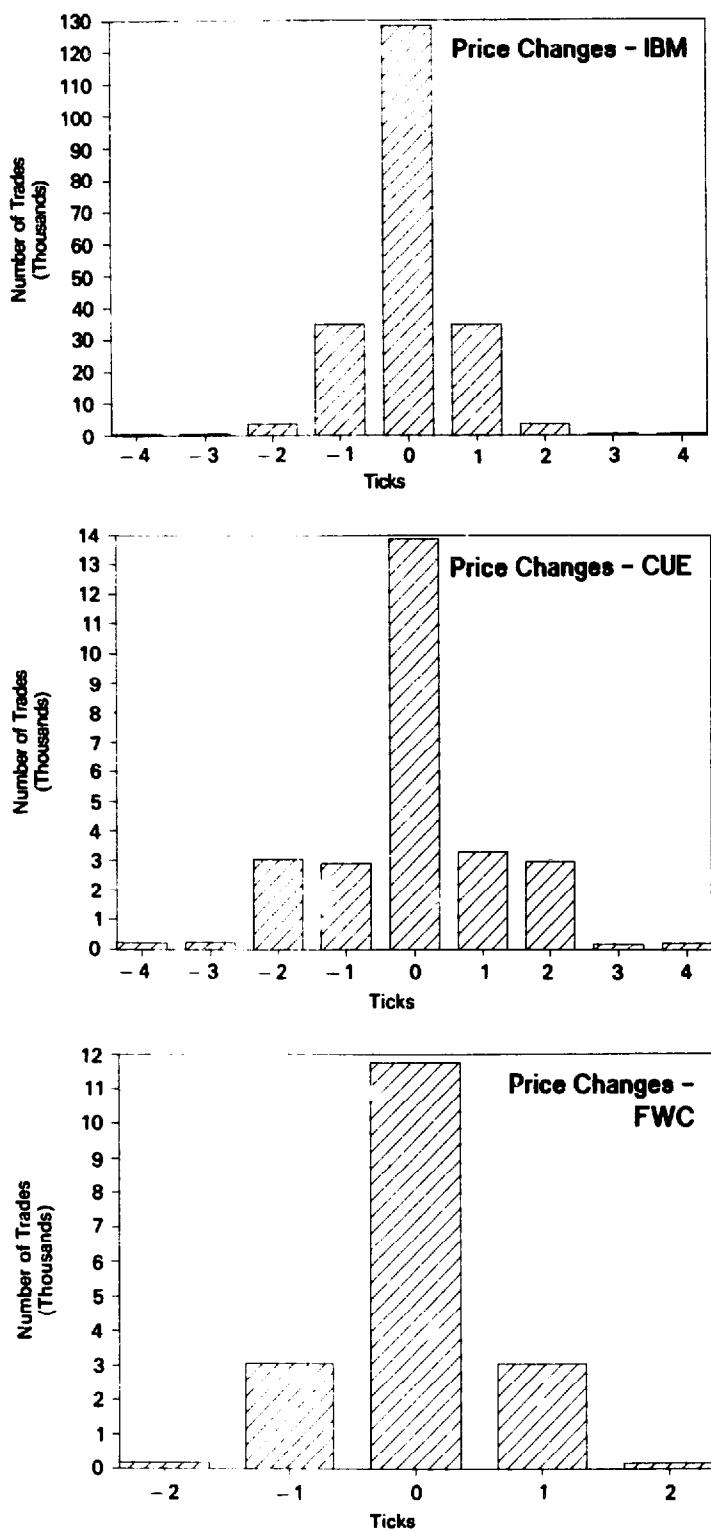


Fig. 2. Histograms of price changes, time-between-trades, and dollar volume of International Business Machines Corporation (IBM - 206,794 trades), Quantum Chemical Corporation (CUE - 26,927 trades), Foster Wheeler Corporation (FWC - 18,199 trades), Handy and Harman Company (HNH - 3,174 trades), Navistar International Corporation (NAV - 96,127 trades), and American Telephone and Telegraph Company (T - 180,726 trades), for the period from 4 January 1988 to 30 December 1988.

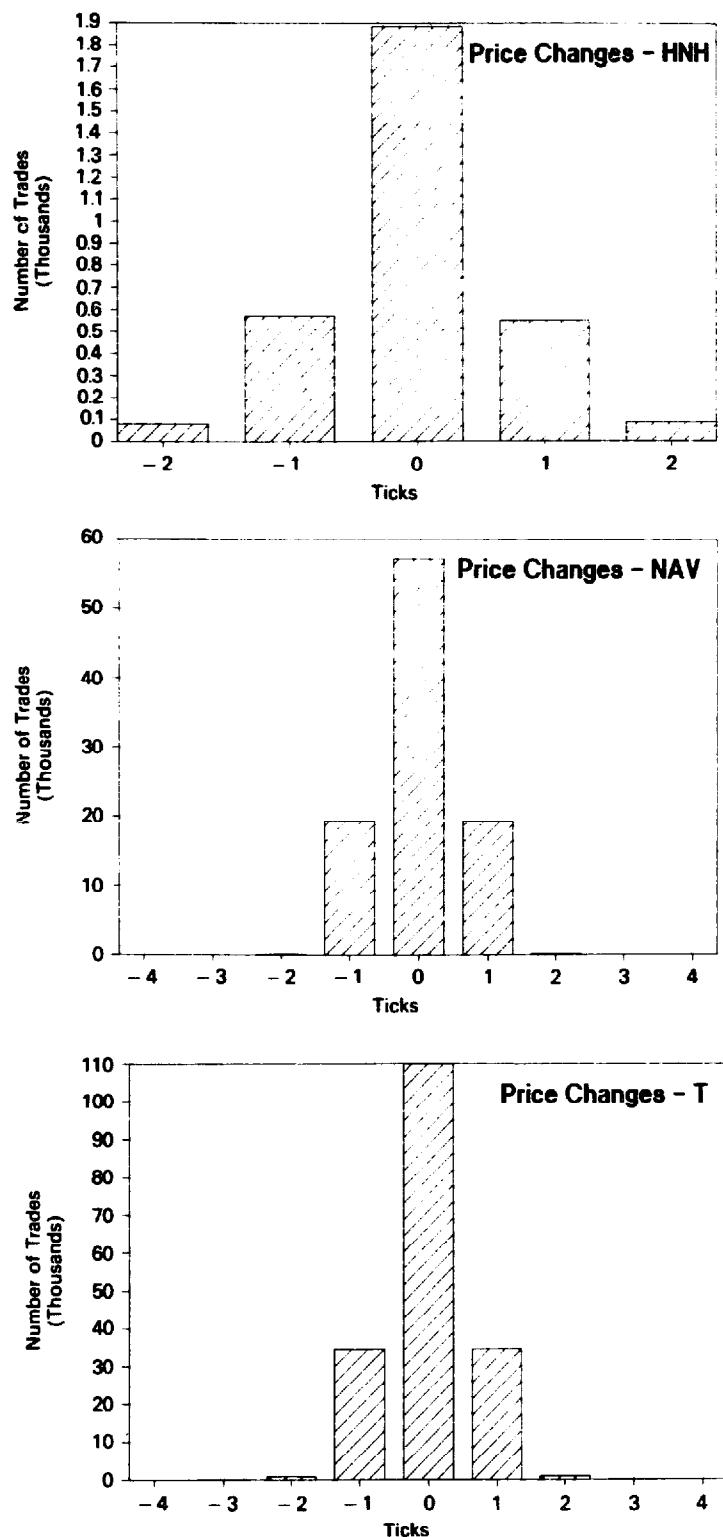


Fig. 2 (continued)

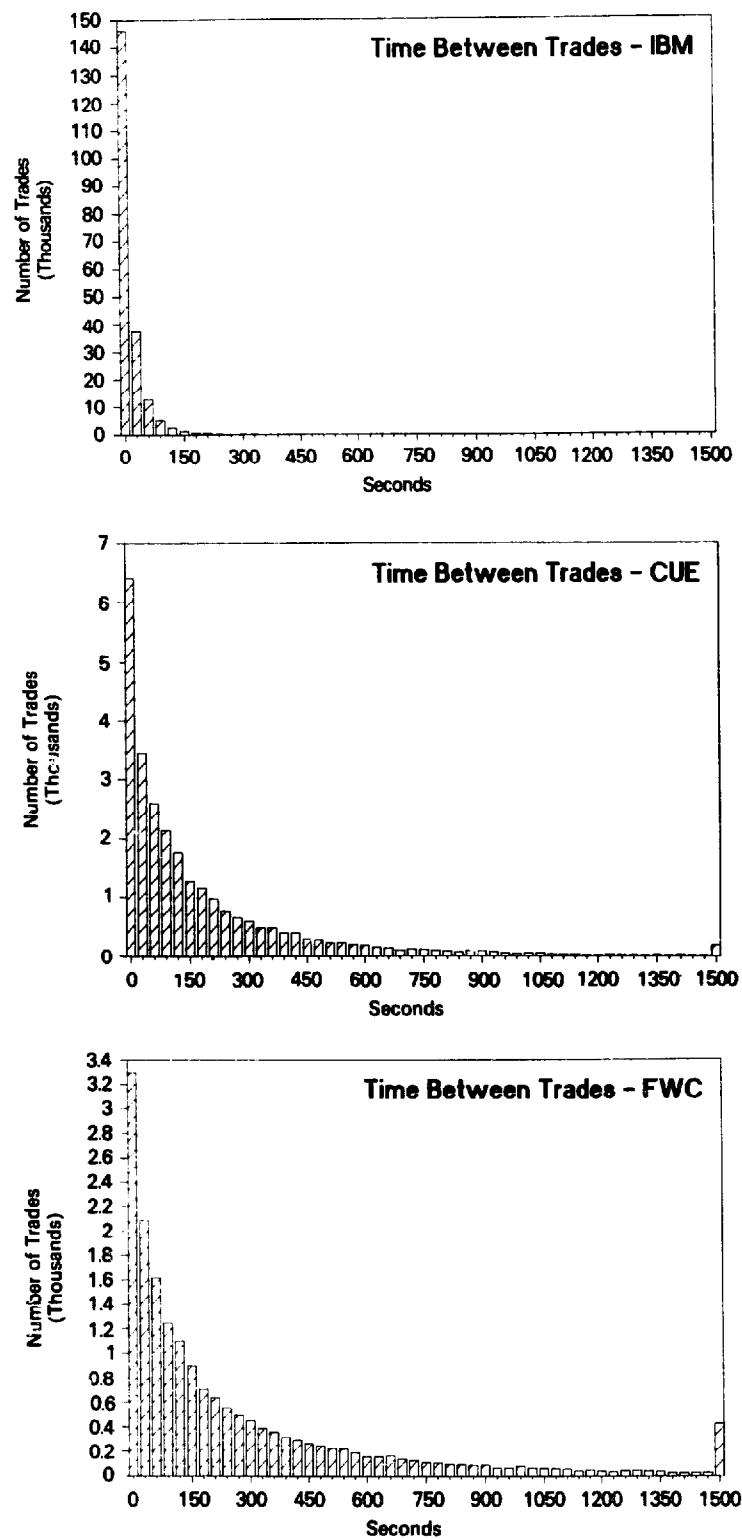


Fig. 2 (continued)

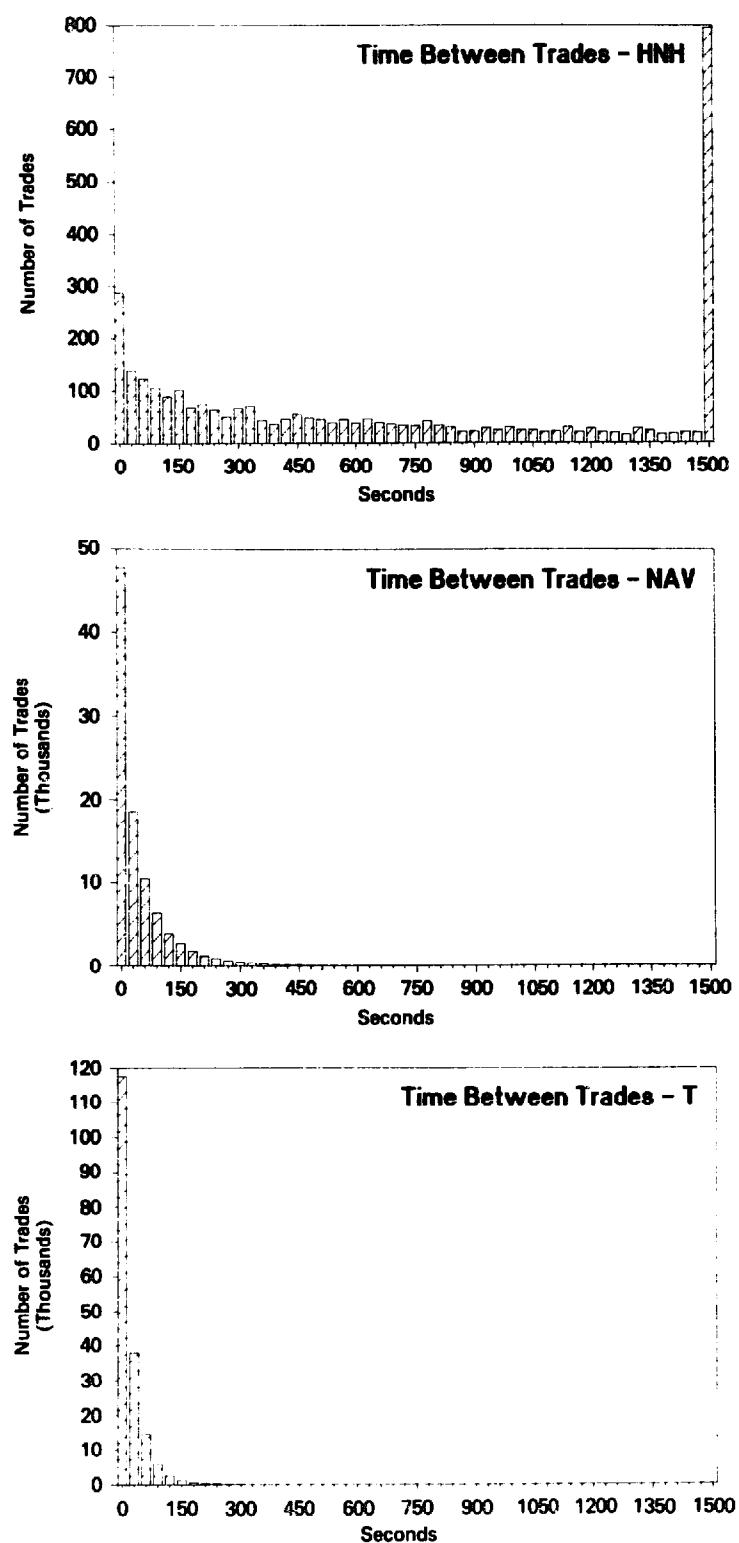


Fig. 2 (continued)

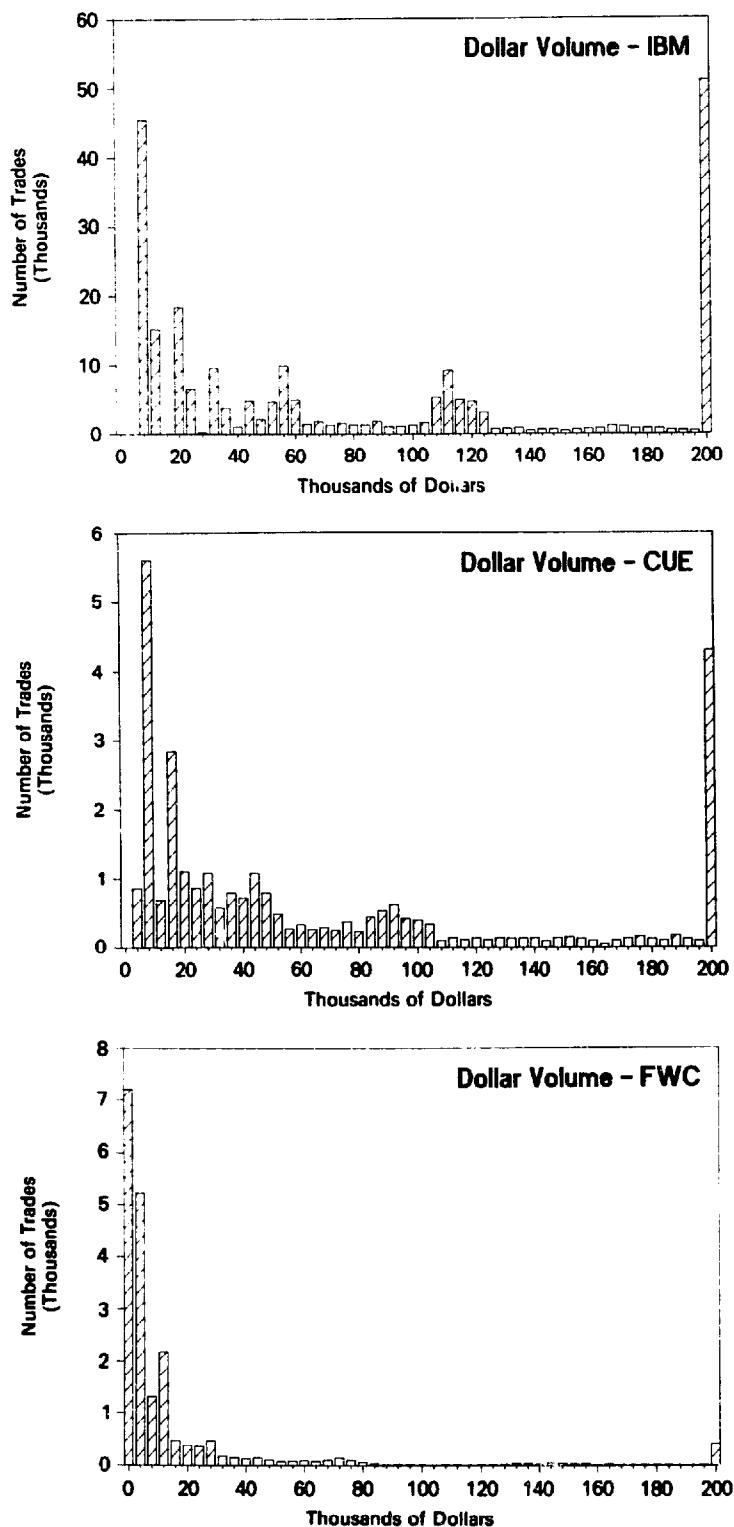


Fig. 2 (continued)

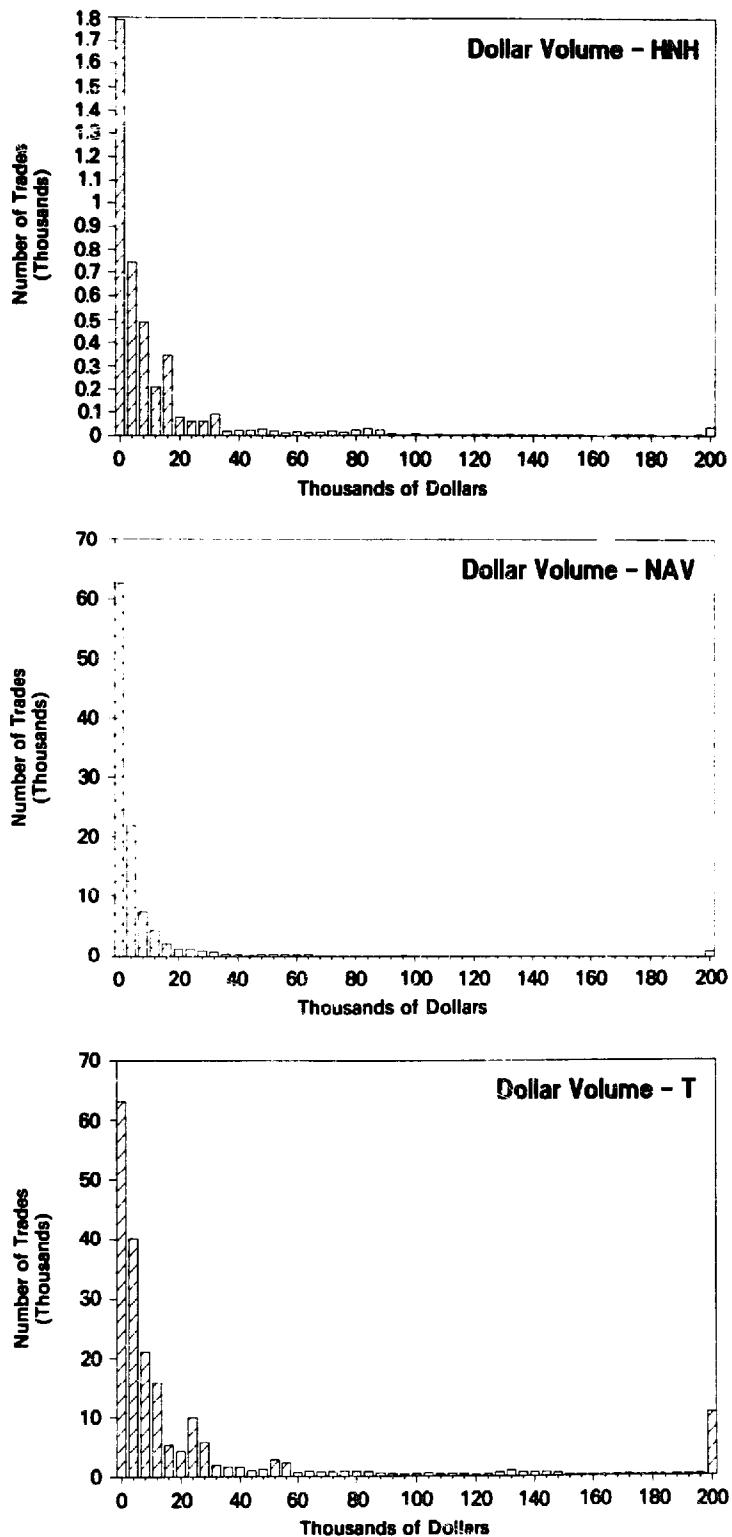


Fig. 2 (continued)

In selecting the explanatory variables  $X_k$ , we seek to capture several aspects of transaction price changes. First, we would like to allow for clock-time effects, since there is currently some dispute over whether trade-to-trade prices are stable in transaction time versus clock time. Second, we would like to account for the effects of the bid/ask spread on price changes, since many transactions are merely movements from the bid price to the ask price or vice versa. If, for example, in a sequence of three trades the first and third were buyer-initiated while the second was seller-initiated, the sequence of transaction prices would exhibit reversals due solely to the bid/ask ‘bounce’. Third, we would like to measure how the conditional distribution of price changes shifts in response to a trade of a given volume, i.e., the price impact per unit volume of trade. And fourth, we would like to capture the effects of ‘systematic’ or market-wide movements in prices on the conditional distribution of an individual stock’s price changes. To address these four issues, we first construct the following variables:

- $\Delta t_k$  = Time elapsed between transactions  $k - 1$  and  $k$ , in seconds.
- $AB_{k-1}$  = Bid/ask spread prevailing at time  $t_{k-1}$ , in ticks.
- $Z_{k-l}$  = Three lags ( $l = 1, 2, 3$ ) of the dependent variable  $Z_k$ . Recall that for  $m = 9$ , price changes less than  $-4$  ticks are set equal to  $-4$  ticks (state  $s_1$ ), and price changes greater than  $+4$  ticks are set equal to  $+4$  ticks (state  $s_9$ ), and similarly for  $m = 5$ .
- $V_{k-l}$  = Three lags ( $l = 1, 2, 3$ ) of the dollar volume of the  $(k - l)$ th transaction, defined as the price of the  $(k - l)$ th transaction (in dollars, not ticks) times the number of shares traded (denominated in 100’s of shares), hence dollar volume is denominated in \$100’s of dollars. To reduce the influence of outliers, if the share volume of a trade exceeds the 99.5 percentile of the empirical distribution of share volume for that stock, we set it equal to the 99.5 percentile.<sup>7</sup>
- $SP500_{k-l}$  = Three lags ( $l = 1, 2, 3$ ) of five-minute continuously-compounded returns of the Standard and Poor’s 500 index futures price, for the contract maturing in the closest month beyond the month in which transaction  $k - l$  occurred, where the return is computed

<sup>7</sup>For example, the 99.5 percentile for IBM’s share volume is 16,500 shares, hence all IBM trades exceeding 16,500 shares are set equal to 16,500 shares. By definition, only one-half of one percent of the 206,794 IBM trades [or 1,034 trades] were ‘censored’ in this manner. We chose not to discard these trades because omitting them could affect our estimates of the lag structure, which is extremely sensitive to the sequence of trades. For the five remaining stocks, the 99.5 percentiles for share volume are: CUE = 21,300, FWC = 31,700, HNH = 20,000, NAV = 50,000, and T = 44,100.

with the futures price recorded one minute before the nearest round minute *prior* to  $t_{k-1}$  and the price recorded five minutes before this. More formally, we have:

$$SP500_{k-1} \equiv \log(F(t_{k-1}^- - 60)/F(t_{k-1}^- - 360)), \quad (4.1)$$

$$SP500_{k-2} \equiv \log(F(t_{k-1}^- - 360)/F(t_{k-1}^- - 660)), \quad (4.2)$$

$$SP500_{k-3} \equiv \log(F(t_{k-1}^- - 660)/F(t_{k-1}^- - 960)), \quad (4.3)$$

where  $F(t^-)$  is the S&P500 index futures price at time  $t^-$  (measured in seconds) for the contract maturing the closest month beyond the month of transaction  $k-l$ , and  $t^-$  is the nearest round minute prior to time  $t$  (for example, if  $t$  is 10:35:47, then  $t^-$  is 10:35:00).<sup>8</sup>

- $IBS_{k-l}$  = Three lags ( $l = 1, 2, 3$ ) of an indicator variable that takes the value 1 if the  $(k-l)$ th transaction price is greater than the average of the quoted bid and ask prices at time  $t_{k-l}$ , the value -1 if the  $(k-l)$ th transaction price is less than the average of the bid and ask prices at time  $t_{k-l}$ , and 0 otherwise, i.e.,

$$IBS_{k-l} \equiv \begin{cases} 1 & \text{if } P_{k-l} > \frac{1}{2}(P_{k-l}^a + P_{k-l}^b), \\ 0 & \text{if } P_{k-l} = \frac{1}{2}(P_{k-l}^a + P_{k-l}^b), \\ -1 & \text{if } P_{k-l} < \frac{1}{2}(P_{k-l}^a + P_{k-l}^b). \end{cases} \quad (4.4)$$

Whether the  $(k-l)$ th transaction price is closer to the ask price or the bid price is one measure of whether the transaction was buyer-initiated ( $IBS_{k-l} = 1$ ) or seller-initiated ( $IBS_{k-l} = -1$ ). If the transaction price is at the midpoint of the bid and ask prices, the indicator is indeterminate ( $IBS_{k-l} = 0$ ).

<sup>8</sup>This rather convoluted timing for computing  $SP500_{k-1}$  ensures that there is no temporal overlap between price changes and the returns to the index futures price. In particular, we first construct a minute-by-minute time series for futures prices by assigning to each round minute the nearest futures transaction price occurring *after* that minute but before the next (hence if the first futures transaction after 10:35:00 occurs at 10:35:15, the futures price assigned to 10:35:00 is this one). If no transaction occurs during this minute, the price prevailing at the previous minute is assigned to the current minute. Then for the price change  $Z_k$ , we compute  $SP500_{k-1}$  using the futures price one minute before the nearest round minute *prior* to  $t_{k-1}$ , and the price five minutes before this (hence if  $t_{k-1}$  is 10:36:45, we use the futures price assigned to 10:35:00 and 10:30:00 to compute  $SP500_{k-1}$ ).

Our specification of  $X'_k \beta$  is then given by the following expression:

$$\begin{aligned}
 X'_k \beta = & \beta_1 \Delta t_k + \beta_2 Z_{k-1} + \beta_3 Z_{k-2} + \beta_4 Z_{k-3} + \beta_5 SP500_{k-1} \\
 & + \beta_6 SP500_{k-2} + \beta_7 SP500_{k-3} + \beta_8 IBS_{k-1} \\
 & + \beta_9 IBS_{k-2} + \beta_{10} IBS_{k-3} + \beta_{11} \{T_\lambda(V_{k-1}) \cdot IBS_{k-1}\} \\
 & + \beta_{12} \{T_\lambda(V_{k-2}) \cdot IBS_{k-2}\} + \beta_{13} \{T_\lambda(V_{k-3}) \cdot IBS_{k-3}\}. \quad (4.5)
 \end{aligned}$$

The variable  $\Delta t_k$  is included in  $X_k$  to allow for clock-time effects on the conditional mean of  $Z_k^*$ . If prices are stable in transaction time rather than clock time, this coefficient should be zero. Lagged price changes are included to account for serial dependencies, and lagged returns of the S&P500 index futures price are included to account for market-wide effects on price changes.

To measure the price impact of a trade per unit volume we include the term  $T_\lambda(V_{k-1})$ , dollar volume transformed according to the Box and Cox (1964) specification  $T_\lambda(\cdot)$ :

$$T_\lambda(x) \equiv (x^\lambda - 1)/\lambda, \quad (4.6)$$

where  $\lambda \in [0, 1]$  is also a parameter to be estimated. The Box–Cox transformation allows dollar volume to enter into the conditional mean nonlinearly, a particularly important innovation since common intuition suggests that price impact may exhibit economies of scale with respect to dollar volume, i.e., although total price impact is likely to increase with volume, the marginal price impact probably does not. The Box–Cox transformation captures the linear specification [ $\lambda = 1$ ] and concave specifications up to and including the logarithmic function [ $\lambda = 0$ ]. The estimated curvature of this transformation will play an important role in the measurement of price impact.

The transformed dollar volume variable is interacted with  $IBS_{k-1}$ , an indicator of whether the trade was buyer-initiated ( $IBS_k = 1$ ), seller-initiated ( $IBS_k = -1$ ), or indeterminate ( $IBS_k = 0$ ). A positive  $\beta_{11}$  would imply that buyer-initiated trades tend to push prices up and seller-initiated trades tend to drive prices down. Such a relation is predicted by several information-based models of trading, e.g., Easley and O’Hara (1987). Moreover, the magnitude of  $\beta_{11}$  is the per-unit volume impact on the conditional mean of  $Z_k^*$ , which may be readily translated into the impact on the conditional probabilities of observed price changes. The sign and magnitudes of  $\beta_{12}$  and  $\beta_{13}$  measure the persistence of price impact.

Finally, to complete our specification we must parametrize the conditional variance  $\sigma_k^2 \equiv \gamma_0^2 + \sum \gamma_i^2 W_{ik}$ . To allow for clock-time effects we include  $\Delta t_k$ , and since there is some evidence linking bid/ask spreads to the information content

and volatility of price changes [see, for example, Glosten (1987), Hasbrouck (1988, 1991a,b), and Petersen and Umlauf (1990)], we also include the lagged spread  $AB_{k-1}$ . Also, recall from section 2.2 that the parameters  $\alpha$ ,  $\beta$ , and  $\gamma$  are unidentified without additional restrictions, hence we make the identification assumption that  $\gamma_0^2 = 1$ . Our variance parametrization is then:

$$\sigma_k^2 \equiv 1 + \gamma_1^2 \Delta t_k + \gamma_2^2 AB_{k-1}. \quad (4.7)$$

In summary, our nine-state specification requires the estimation of 24 parameters: the partition boundaries  $\alpha_1, \dots, \alpha_8$ , the variance parameters  $\gamma_1$  and  $\gamma_2$ , the coefficients of the explanatory variables  $\beta_1, \dots, \beta_{13}$ , and the Box-Cox parameter  $\lambda$ . The five-state specification requires the estimation of only 20 parameters.

## 5. The maximum likelihood estimates

We compute the maximum likelihood estimators numerically using the algorithm proposed by Berndt, Hall, Hall, and Hausman (1974), hereafter BHHH. The advantage of BHHH over other search algorithms is its reliance on only first derivatives, an important computational consideration for sample sizes such as ours.

The asymptotic covariance matrix of the parameter estimates was computed as the negative inverse of the matrix of [numerically determined] second derivatives of the log-likelihood function with respect to the parameters, evaluated at the maximum likelihood estimates. We used a tolerance of 0.001 for the convergence criterion suggested by BHHH [the product of the gradient and the direction vector]. To check the robustness of our numerical search procedure, we used several different sets of starting values for each stock, and in all instances our algorithm converged to virtually identical parameter estimates.

All computations were performed in double precision in an ULTRIX environment on a DEC 5000/200 workstation with 16 Mb of memory, using our own FORTRAN implementation of the BHHH algorithm with analytical first derivatives. As a rough guide to the computational demands of ordered probit, note that the numerical estimation procedure for the stock with the largest number of trades (IBM, with 206,794 trades) required only 2 hours and 45 minutes of cpu time.

In table 2a, we report the maximum likelihood estimates of the ordered probit model for our six stocks. Entries in each of the columns labeled with ticker symbols are the parameter estimates for that stock, and to the immediate right of each parameter estimate is the corresponding z-statistic, which is asymptotically distributed as a standard normal variate under the null hypothesis that the coefficient is zero, i.e., it is the parameter estimate divided by its asymptotic standard error.

Table 2a shows that the partition boundaries are estimated with high precision for all stocks. As expected, the  $z$ -statistics are much larger for those stocks with many more observations. The parameters for  $\sigma_k^2$  are also statistically significant, hence homoskedasticity may be rejected at conventional significance levels; larger bid/ask spreads and longer time intervals increase the conditional volatility of the disturbance.

The conditional means of the  $Z_k^*$ 's for all stocks are only marginally affected by  $\Delta t_k$ . Moreover, the  $z$ -statistics are minuscule, especially in light of the large sample sizes. However, as mentioned above,  $\Delta t$  does enter into the  $\sigma_k^2$  expression significantly, hence clock time is important for the conditional variances, but not for the conditional means of  $Z_k^*$ . Note that this does not necessarily imply the same for the conditional distribution of the  $Z_k$ 's, which is *nonlinearly* related to the conditional distribution of the  $Z_k^*$ 's. For example, the conditional mean of the  $Z_k$ 's may well depend on the conditional variance of the  $Z_k^*$ 's, so that clock time can still affect the conditional mean of observed price changes even though it does not affect the conditional mean of  $Z_k^*$ .

More striking is the significance and sign of the lagged price change coefficients  $\hat{\beta}_2$ ,  $\hat{\beta}_3$ , and  $\hat{\beta}_4$ , which are negative for all stocks, implying a tendency towards price reversals. For example, if the past three price changes were each one tick, the conditional mean of  $Z_k^*$  changes by  $\hat{\beta}_2 + \hat{\beta}_3 + \hat{\beta}_4$ . However, if the sequence of price changes was  $1/ - 1/1$ , then the effect on the conditional mean is  $\hat{\beta}_2 - \hat{\beta}_3 + \hat{\beta}_4$ , a quantity closer to zero for each of the security's parameter estimates.<sup>9</sup>

Note that these coefficients measure reversal tendencies beyond that induced by the presence of a constant bid/ask spread as in Roll (1984). The effect of this 'bid/ask bounce' on the conditional mean should be captured by the indicator variables  $IBS_{k-1}$ ,  $IBS_{k-2}$ , and  $IBS_{k-3}$ . In the absence of all other information [such as market movements, past price changes, etc.], these variables pick up any price effects that buys and sells might have on the conditional mean. As expected, the estimated coefficients are generally negative, indicating the presence of reversals due to movements from bid to ask or ask to bid prices. In section 6.1 we shall compare their magnitudes explicitly, and conclude that the conditional mean of price changes is *path-dependent* with respect to past price changes.

The lagged S&P500 returns are also significant, but have a more persistent effect on some securities. For example, the coefficient for the first lag of the S&P500 is large and significant for IBM, but the coefficient for the third

<sup>9</sup>In an earlier specification, in place of lagged price changes we included separate indicator variables for eight of the nine states of each lagged price change. But because the coefficients of the indicator variables increased monotonically from the  $-4$  state to the  $+4$  state [state 0 was omitted] in almost exact proportion to the tick change, we chose the more parsimonious specification of including the actual lagged price change.

Table 2a

Maximum likelihood estimates of the ordered probit model for transaction price changes of International Business Machines Corporation (IBM – 206,794 trades), Quantum Chemical Corporation (CUE – 26,927 trades), Foster Wheeler Corporation (FWC – 18,199 trades), Handy and Harman Company (HNH – 3,174 trades), Navistar International Corporation (NAV – 96,127 trades), and American Telephone and Telegraph Company (T – 180,726 trades), for the period from 4 January 1988 to 30 December 1988. Each z-statistic is asymptotically standard normal under the null hypothesis that the corresponding coefficient is zero.

Parameter	IBM	z	CUE	z	FWC <sup>b</sup>	z	HNH <sup>b</sup>	z	NAV	z	T	z
Partition boundaries <sup>a</sup>												
$\alpha_1$	-4.670	-145.65	-6.213	-18.92	-4.378	-25.24	-4.456	-5.98	-7.263	-39.23	-8.073	-56.95
$\alpha_2$	-4.157	-157.75	-5.447	-18.99	-1.712	-25.96	-1.801	-5.92	-7.010	-36.53	-7.270	-62.40
$\alpha_3$	-3.109	-171.59	-2.795	-19.14	1.679	26.32	1.923	5.97	-6.251	-37.22	-5.472	-63.43
$\alpha_4$	-1.344	-155.47	-1.764	-18.95	4.334	25.26	4.477	5.85	-1.972	-34.59	-1.850	-61.41
$\alpha_5$	1.326	154.91	1.605	18.81	-	-	-	-	1.938	34.66	1.977	62.82
$\alpha_6$	3.126	167.81	2.774	19.11	-	-	-	-	6.301	36.36	5.378	62.43
$\alpha_7$	4.205	152.17	5.502	19.10	-	-	-	-	7.742	31.63	7.294	57.63
$\alpha_8$	4.732	138.75	6.150	18.94	-	-	-	-	8.638	30.26	8.156	56.23
$\gamma_1 : \Delta I/100$	0.399	15.57	0.499	11.62	0.275	11.26	0.187	4.07	0.428	10.01	0.387	8.89
$\gamma_2 : AB^{-1}$	0.515	71.08	1.110	15.39	0.723	14.54	1.109	4.48	0.869	19.93	0.868	38.16
$\beta_1 : \Delta I/100$	-0.115	-11.42	-0.014	-2.14	-0.013	-3.50	-0.010	-2.69	-0.032	-3.82	-0.127	-9.51
$\beta_2 : Z^{-1}$	-1.012	-135.57	-0.333	-13.46	-1.325	-24.49	-0.740	-5.18	-2.609	-36.32	-2.346	-62.74
$\beta_3 : Z^{-2}$	-0.532	-85.00	-0.000	-0.03	-0.638	-16.45	-0.406	-4.06	-1.521	-34.13	-1.412	-56.52
$\beta_4 : Z^{-3}$	-0.211	-47.15	-0.020	-1.42	-0.223	-9.23	-0.116	-1.84	-0.536	-31.63	-0.501	-47.91
$\beta_5 : SP500_{-1}$	1.120	54.22	2.292	1.54	1.359	13.49	0.472	1.36	0.419	8.05	0.625	17.12
$\beta_6 : SP500_{-2}$	-0.257	-12.06	1.373	9.61	0.302	2.93	0.448	1.20	0.150	2.87	0.177	4.96
$\beta_7 : SP500_{-3}$	0.006	-0.26	0.677	5.15	0.204	1.97	0.388	1.13	0.159	3.02	0.141	3.93
$\beta_8 : IB\$^{-1}$	-1.137	-63.64	-1.915	-15.36	-0.791	-7.81	-0.803	-2.89	-0.501	-17.38	-0.740	-23.01
$\beta_9 : IB\$^{-2}$	-0.369	-21.55	-0.279	-3.37	-0.184	-3.66	-0.184	-0.75	-0.370	-15.38	-0.340	-18.11
$\beta_{10} : IB\$^{-3}$	-0.174	-10.29	0.079	0.98	-0.177	-3.64	-0.022	-0.17	-0.301	-15.37	-0.299	-19.78
$\beta_{11} : T_\lambda(V_{-1})^1IB\$^{-1}$	0.122	47.37	0.217	12.97	0.050	1.80	0.038	0.55	0.013	2.56	0.032	4.51
$\beta_{12} : T_\lambda(V_{-2})^1IB\$^{-2}$	0.047	18.57	0.036	2.83	0.015	1.54	0.036	0.55	0.011	2.54	0.014	4.22
$\beta_{13} : T_\lambda(V_{-3})^1IB\$^{-3}$	0.019	7.70	0.007	0.59	0.015	1.56	-0.006	-0.34	0.005	2.09	0.005	3.02
$\lambda$	0	-	0	-	0.165	1.58	0.191	0.55	0.277	3.50	0.182	5.00

<sup>a</sup>According to the ordered probit model, if the 'virtual' price change  $Z_k^*$  is less than  $\alpha_1$ , then the observed price change is -4 ticks or less; if  $Z_k^*$  is between  $\alpha_1$  and  $\alpha_2$ , then the observed price change is -3 ticks; and so on.

<sup>b</sup>The ordered probit specification for FWC and HNH contains only five states (-2 ticks or less, -1, 0, +1, +2 ticks or more), hence only four  $\alpha$ 's were required. Box-Cox transformation of lagged dollar volume multiplied by the lagged buy/sell indicator, where the Box-Cox parameter  $\lambda$  is estimated jointly with the other ordered probit parameters via maximum likelihood. The Box-Cox parameter  $\lambda$  determines the degree of curvature that the transformation  $T_\lambda(\cdot)$  exhibits in transforming dollar volume  $V_k$  before inclusion as an explanatory variable in the ordered probit specification. If  $\lambda = 1$ , the transformation  $T_\lambda(\cdot)$  is linear, hence dollar volume enters the ordered probit model linearly. If  $\lambda = 0$ , the transformation is equivalent to  $\log(\cdot)$ , hence the natural logarithm of dollar volume enters the ordered probit model. When  $\lambda$  is between 0 and 1, the curvature of  $T_\lambda(\cdot)$  is between logarithmic and linear.

is small and insignificant. However, for the less actively traded stocks such as CUE, all three coefficients are significant and are about the same order of magnitude. As a measure of how quickly market-wide information is impounded into prices, these coefficients confirm the common intuition that smaller stocks react more slowly than larger stocks, which is consistent with the lead/lag effects uncovered by Lo and MacKinlay (1990a).

### 5.1. Diagnostics

A common diagnostic for the specification of an ordinary least squares regression is to examine the properties of the residuals. If, for example, a time series regression is well-specified, the residuals should approximate white noise and exhibit little serial correlation. In the case of ordered probit, we cannot calculate the residuals directly since we cannot observe the latent dependent variable  $Z_k^*$  and therefore cannot compute  $Z_k^* - X_k' \hat{\beta}$ . However, we do have an estimate of the conditional distribution of  $Z_k^*$ , conditioned on the  $X_k$ 's, based on the ordered probit specification and the maximum likelihood parameter estimates. From this we can obtain an estimate of the conditional distribution of the  $\varepsilon_k$ 's, from which we can construct *generalized residuals*  $\hat{\varepsilon}_k$  along the lines suggested by Gourieroux, Monfort, and Trognon (1985):

$$\hat{\varepsilon}_k \equiv E[\varepsilon_k | Z_k, X_k, W_k; \hat{\theta}_{ml}], \quad (5.1)$$

where  $\hat{\theta}_{ml}$  is the maximum likelihood estimator of the unknown parameter vector containing  $\hat{\alpha}$ ,  $\hat{\gamma}$ ,  $\hat{\beta}$ , and  $\hat{\lambda}$ . In the case of ordered probit, if  $Z_k$  is in the  $j$ th state, i.e.,  $Z_k = s_j$ , then the generalized residual  $\hat{\varepsilon}_k$  may be expressed explicitly using the moments of the truncated normal distribution as

$$\begin{aligned} \hat{\varepsilon}_k &= E[\varepsilon_k | Z_k = s_j, X_k, W_k; \hat{\theta}_{ml}] \\ &= \hat{\sigma}_k \cdot \frac{\phi(c_1) - \phi(c_2)}{\Phi(c_2) - \Phi(c_1)}, \end{aligned} \quad (5.2)$$

$$c_1 \equiv \frac{1}{\hat{\sigma}_k} (\hat{\alpha}_{j-1} - X_k' \hat{\beta}), \quad (5.3)$$

$$c_2 \equiv \frac{1}{\hat{\sigma}_k} (\hat{\alpha}_j - X_k' \hat{\beta}), \quad (5.4)$$

$$\hat{\sigma}_k \equiv \sqrt{1 + \hat{\gamma}_1^2 \Delta t_k + \hat{\gamma}_2^2 A B_{k-1}}, \quad (5.5)$$

where  $\phi(\cdot)$  is the standard normal probability density function and for notational convenience, we define  $\alpha_0 \equiv -\infty$  and  $\alpha_m \equiv +\infty$ . Gourieroux et al. (1985) show that these generalized residuals may be used to test for misspecification in a variety of ways. However, some care is required in performing such tests. For example, although a natural statistic to calculate is the first-order autocorrelation of the  $\hat{\varepsilon}_k$ 's, Gourieroux et al. observe that the theoretical autocorrelation of the generalized residuals does not in general equal the theoretical autocorrelation of the  $\varepsilon_k$ 's. Moreover, if the source of serial correlation is an omitted lagged endogenous variable [if, for example, we included too few lags of  $Z_k$  in  $X_k$ ], then further refinements of the usual specification tests are necessary.

Gourieroux et al. derive valid tests for serial correlation from lagged endogenous variables using the *score statistic*, essentially the derivative of the likelihood function with respect to an autocorrelation parameter, evaluated at the maximum likelihood estimates under the null hypothesis of no serial correlation. More specifically, consider the following model for our  $Z_k^*$ :

$$Z_k^* = \varphi Z_{k-1}^* + X_k' \beta + \varepsilon_k, \quad |\varphi| < 1. \quad (5.6)$$

In this case, the score statistic  $\hat{\xi}_1$  is the derivative of the likelihood function with respect to  $\varphi$  evaluated at the maximum likelihood estimates. Under the null hypothesis that  $\varphi = 0$ , it simplifies to the following expression:

$$\hat{\xi}_1 \equiv \left( \sum_{k=2}^n \hat{Z}_{k-1} \hat{\varepsilon}_k \right)^2 / \sum_{k=2}^n \hat{Z}_{k-1}^2 \hat{\varepsilon}_k^2, \quad (5.7)$$

where

$$\hat{Z}_k \equiv E[Z_k^* | Z_k, X_k, W_k; \hat{\theta}_{ml}] \quad (5.8)$$

$$= X_k' \hat{\beta} + \hat{\varepsilon}_k. \quad (5.9)$$

When  $\varphi = 0$ ,  $\hat{\xi}_1$  is asymptotically distributed as a  $\chi_1^2$  variate. Therefore, using  $\hat{\xi}_1$  we can test for the presence of autocorrelation induced by the omitted variable  $Z_{k-1}^*$ . More generally, we can test the higher-order specification:

$$Z_k^* = \varphi Z_{k-j}^* + X_k' \beta + \varepsilon_k, \quad |\varphi| < 1, \quad (5.10)$$

by using the score statistic  $\hat{\xi}_j$ ,

$$\hat{\xi}_j \equiv \left( \sum_{k=j+1}^n \hat{Z}_{k-j} \hat{\varepsilon}_k \right)^2 / \sum_{k=j+1}^n \hat{Z}_{k-j}^2 \hat{\varepsilon}_k^2, \quad (5.11)$$

which is also asymptotically  $\chi_1^2$  under the null hypothesis that  $\varphi = 0$ .

For further intuition, we can compute the sample correlation  $\hat{v}_j$  of the generalized residual  $\hat{\varepsilon}_k$  with the lagged generalized fitted values  $\hat{Z}_{k-j}$ . Under the null hypothesis of no serial correlation in the  $\varepsilon_k$ 's, the theoretical value of this correlation is zero, hence the sample correlation will provide one measure of the economic impact of misspecification. These are reported in table 2b for our sample of six stocks, and they are all quite small, ranging from  $-0.088$  to  $0.030$ .

Finally, table 2c reports the score statistics  $\xi_j, j = 1, \dots, 12$ . Since we have included three lags of  $Z_k$  in our specification of  $X_k$ , it is no surprise that none of the score statistics for  $j = 1, 2, 3$  are statistically significant at the 5% level. However, at lag 4, the score statistics for all stocks except CUE and HNH are significant, indicating the presence of some serial dependence not accounted for by our specification. But recall that we have very large sample sizes so that virtually any point null hypothesis will be rejected. With this in mind, the score statistics seem to indicate a reasonably good fit for all but one stock, NAV, whose score statistic is significant at every lag, suggesting the need for respecification. Turning back to the cross-autocorrelations reported in table 2b, we see that NAV's residual  $\hat{\varepsilon}_k$  has a  $-0.088$  correlation with  $\hat{Z}_{k-4}$ , the largest in table 2b in absolute value. This suggests that adding  $Z_{k-4}$  as a regressor might improve the specification for NAV.

There are a number of other specification tests that can check the robustness of the ordered probit specification, but they should be performed with an eye towards particular applications. For example, when studying the impact of information variables on volatility, a more pressing concern would be the specification of the conditional variance  $\sigma_k^2$ . If some of the parameters have important economic interpretations, their stability can be checked by simple likelihood ratio tests on subsamples of the data. If forecasting price changes is of interest, an  $R^2$ -like measure can readily be constructed to measure how much variability can be explained by the predictors. The ordered probit model is flexible enough to accommodate virtually any specification test designed for simple regression models, but has many obvious advantages over OLS as we shall see below.

### *5.2. Endogeneity of $\Delta t_k$ and $IBS_k$*

Our inferences in the preceding sections are based on the implicit assumption that the explanatory variables  $X_k$  are all exogenous or predetermined with respect to the dependent variable  $Z_k$ . However, the variable  $\Delta t_k$  is contemporaneous to  $Z_k$  and deserves further discussion.

Recall that  $Z_k$  is the price change between trades at time  $t_{k-1}$  and time  $t_k$ . Since  $\Delta t_k$  is simply  $t_k - t_{k-1}$ , it may well be that  $\Delta t_k$  and  $Z_k$  are determined simultaneously, in which case our parameter estimates are generally inconsistent. In fact, there are several plausible arguments for the endogeneity of  $\Delta t_k$  [see,

Table 2b  
 Cross-autocorrelation coefficients  $\hat{v}_j$ ,  $j = 1, \dots, 12$ , of generalized residuals  $\{\hat{e}_k\}$  with lagged generalized fitted price changes  $\hat{Z}_{k-j}$  from the ordered probit estimation for transaction price changes of International Business Machines Corporation (IBM - 206,794 trades), Quantum Chemical Corporation (CUE - 26,927 trades), Foster Wheeler Corporation (FWC - 18,199 trades), Handy and Harman Company (HNH - 3,174 trades), Navistar International Corporation (NAV - 96,127 trades), and American Telephone and Telegraph Company (T - 180,726 trades), for the period from 4 January 1988 to 30 December 1988.<sup>a</sup>

Stock	$\hat{v}_1$	$\hat{v}_2$	$\hat{v}_3$	$\hat{v}_4$	$\hat{v}_5$	$\hat{v}_6$	$\hat{v}_7$	$\hat{v}_8$	$\hat{v}_9$	$\hat{v}_{10}$	$\hat{v}_{11}$	$\hat{v}_{12}$
IBM	-0.005	0.002	0.005	-0.043	-0.008	0.001	-0.001	0.001	0.000	-0.001	-0.005	0.000
CUE	-0.008	0.001	-0.006	0.010	0.013	0.003	0.006	0.008	-0.002	0.004	-0.004	0.000
FWC	-0.006	0.000	0.007	-0.032	-0.001	-0.007	-0.004	-0.004	-0.003	-0.003	0.013	-0.004
HNH	-0.012	-0.007	0.007	-0.027	-0.009	0.012	0.019	-0.001	0.009	0.030	-0.018	0.018
NAV	0.005	0.014	0.020	-0.088	-0.011	-0.014	-0.016	-0.011	-0.010	-0.013	-0.009	-0.014
T	0.002	0.013	0.015	-0.080	-0.005	-0.011	-0.006	-0.007	-0.007	-0.006	-0.001	-0.006

<sup>a</sup> If the ordered probit model is correctly specified, these cross-autocorrelations should be close to zero.

Table 2c

Score test statistics  $\hat{\xi}_j$ ,  $j = 1, \dots, 12$ , where  $\hat{\xi}_j \stackrel{a}{\sim} \chi_1^2$  under the null hypothesis of no serial correlation in the ordered probit disturbances  $\{\hat{\varepsilon}_k\}$ , using the generalized residuals  $\{\hat{\varepsilon}_k\}$  from ordered probit estimation for transaction price changes of International Business Machines Corporation (IBM – 206,794 trades), Quantum Chemical Corporation (CUE – 26,927 trades), Foster Wheeler Corporation (FWC – 18,199 trades), Handy and Harman Company (HNH – 3,174 trades), Navistar International Corporation (NAV – 96,127 trades), and American Telephone and Telegraph Company (T – 180,726 trades), for the period from 4 January 1988 to 30 December 1988.<sup>a</sup>

Stock	$\hat{\xi}_1$ (p-value)	$\hat{\xi}_2$ (p-value)	$\hat{\xi}_3$ (p-value)	$\hat{\xi}_4$ (p-value)	$\hat{\xi}_5$ (p-value)	$\hat{\xi}_6$ (p-value)	$\hat{\xi}_7$ (p-value)	$\hat{\xi}_8$ (p-value)	$\hat{\xi}_9$ (p-value)	$\hat{\xi}_{10}$ (p-value)	$\hat{\xi}_{11}$ (p-value)	$\hat{\xi}_{12}$ (p-value)
IBM	3.29 (0.07)	0.94 (0.33)	3.40 (0.07)	313.12 (0.00)	9.71 (0.66)	0.19 (0.60)	0.28 (0.62)	0.25 (1.00)	0.00 (0.65)	0.21 (1.00)	3.76 (0.05)	0.03 (0.86)
CUE	1.25 (0.26)	0.01 (0.92)	0.72 (0.40)	2.39 (0.12)	4.01 (0.05)	0.24 (0.62)	0.94 (0.33)	1.54 (0.21)	0.11 (0.74)	0.41 (0.52)	0.32 (0.57)	0.00 (1.00)
FWC	0.58 (0.45)	0.00 (1.00)	0.82 (0.37)	17.42 (0.00)	0.02 (0.89)	0.75 (0.39)	0.21 (0.65)	0.14 (0.71)	0.15 (0.70)	0.16 (0.69)	3.01 (0.08)	0.27 (0.60)
HNH	0.35 (0.55)	0.13 (0.72)	0.15 (0.70)	2.10 (0.15)	0.22 (0.64)	0.40 (0.53)	1.06 (0.30)	0.00 (1.00)	0.24 (0.62)	2.60 (0.11)	0.96 (0.33)	1.00 (0.32)
NAV	2.37 (0.12)	17.50 (0.00)	38.00 (0.00)	684.06 (0.00)	11.76 (0.00)	18.20 (0.00)	22.38 (0.00)	11.72 (0.00)	9.95 (0.00)	14.61 (0.00)	7.14 (0.01)	17.02 (0.00)
T	0.94 (0.33)	30.12 (0.00)	40.42 (0.00)	1003.69 (0.00)	3.02 (0.08)	17.87 (0.00)	4.96 (0.03)	6.22 (0.01)	7.29 (0.01)	5.52 (0.02)	0.04 (0.84)	4.89 (0.03)

<sup>a</sup> If the ordered probit model is correctly specified, these test statistics should follow a  $\chi_1^2$  statistic which falls in the interval [0.00, 3.84] with 95% probability.

for example, Admati and Pfleiderer (1988, 1989) and Easley and O'Hara (1990)]. One such argument turns on the tendency of floor brokers to break up large trades into smaller ones, and time the executions carefully during the course of the day or several days. By 'working' the order, the floor broker can minimize the price impact of his trades and obtain more favorable execution prices for his clients. But by selecting the times between his trades based on current market conditions, which include information also affecting price changes, the floor broker is creating endogenous trade times.

However, any given sequence of trades in our dataset does not necessarily correspond to consecutive transactions of any single individual [other than the specialist of course], but is the result of many buyers and sellers interacting with the specialist. For example, even if a floor broker were working a larger order, in between his orders might be purchases and sales from other floor brokers, market orders, and triggered limit orders. Therefore, the  $\Delta t_k$ 's also reflect these trades, which are not necessarily information-motivated.

Another more intriguing reason that  $\Delta t_k$  may be exogenous is that floor brokers have an economic incentive to minimize the correlation between  $\Delta t_k$  and virtually all other exogenous and predetermined variables. To see this, suppose the floor broker timed his trades in response to some exogenous variable also affecting price changes, call it 'weather'. Suppose that price changes tend to be positive in good weather and negative in bad weather. Knowing this, the floor broker will wait until bad weather prevails before buying, hence trade times and price changes are simultaneously determined by weather. However, if other traders are also aware of these relations, they can garner information about the floor broker's intent by watching his trades and by recording the weather, and trade against him successfully. To prevent this, the floor broker must trade to deliberately minimize the correlation between his trade times and the weather. Therefore, the floor broker has an economic incentive to reduce simultaneous equations bias! Moreover, this argument applies to any other economic variable that can be used to jointly forecast trade times and price changes. For these two reasons, we assume that  $\Delta t_k$  is exogenous.

We have also explored some adjustments for the endogeneity of  $\Delta t_k$  along the lines of Hausman (1978) and Newey (1985), and our preliminary estimates show that although exogeneity of  $\Delta t_k$  may be rejected at conventional significance levels [recall our sample sizes], the estimates do not change much once endogeneity is accounted for by an instrumental variables estimation procedure.

There are, however, other contemporaneous variables that we would like to include as regressors which cannot be deemed exogenous (see the discussion of  $IBS_k$  in section 6.2 below), and for these we must wait until the appropriate econometric tools become available.

## 6. Applications

In applying the ordered probit model to particular issues of the market microstructure, we must first consider how to interpret its parameter estimates from an economic perspective. Since ordered probit may be viewed as a generalization of a linear regression model to situations with a discrete dependent variable, interpreting its parameter estimates is much like interpreting coefficients of a linear regression: the particular interpretation depends critically on the underlying economic motivation for including and excluding the specific regressors.

In a very few instances, theoretical paradigms might yield testable implications in the form of linear regression equations, e.g., the CAPM's security market line. In most cases, however, linear regression is used to capture and summarize empirical relations in the data that have not yet been derived from economic first principles. In much the same way, ordered probit may be interpreted as a means of capturing and summarizing relations among price changes and other economic variables such as volume. Such relations have been derived from first principles only in the most simplistic and stylized of contexts, under very specific and, therefore, often counterfactual assumptions about agents' preferences, information sets, alternative investment possibilities, sources of uncertainty and their parametric form (usually Gaussian), and the timing and allowable volume and type of trades.<sup>10</sup> Although such models do yield valuable insights about the economics of the market microstructure, they are too easily rejected by the data because of the many restrictive assumptions needed to obtain readily interpretable closed-form results.

Nevertheless, the broader implications of such models can still be 'tested' by checking for simple relations among economic quantities, as we illustrate in section 6.1. However, some care must be taken in interpreting such results, as in the case of a simple linear regression of prices on quantities which cannot be interpreted as an estimated demand curve without imposing additional economic structure.

In particular, although the ordered probit model can shed light on how price changes respond to specific economic variables, it cannot give us economic insights beyond whatever structure we choose to impose *a priori*. For example, since we have placed no specific theoretical structure on how prices are formed, our ordered probit estimates cannot yield sharp implications for the impact of floor brokers 'working' an order [executing a large order in smaller bundles to obtain the best average price]. The ordered probit estimates will reflect the combined actions and interactions of these floor brokers, the specialists, and

<sup>10</sup>Just a few recent examples of this growing literature are Amihud and Mendelson (1980), Admati and Pfleiderer (1988, 1989), Easley and O'Hara (1987), Garman (1976), Glosten and Milgrom (1985), Grundy and McNichols (1989), Ho and Stoll (1980, 1981), Karpoff (1986), Kyle (1985), Stoll (1989), and Wang (1992).

individual and institutional investors, all trading with and against each other. Unless we are estimating a fully articulated model of economic equilibrium that contains these kinds of market participants, we cannot separate their individual impact in determining price changes. For example, without additional structure we cannot answer the question: What is the price impact of an order that is *not* 'worked'?

However, if we were able to identify those large trades that did benefit from the services of a floor broker, we could certainly compare and contrast their empirical price dynamics with those of 'unworked' trades using the ordered probit model. Such comparisons might provide additional guidelines and restrictions for developing new theories of the market microstructure. Interpreted in this way, the ordered probit model can be a valuable tool for uncovering empirical relations even in the absence of a highly parametrized theory of the market microstructure. To illustrate this aspect of ordered probit, in the following section we consider three specific applications of the parameter estimates of section 5: a test for order-flow dependence in price changes, a measure of price impact, and a comparison of ordered probit to ordinary least squares.

### *6.1. Order-flow dependence*

Several recent theoretical papers in the market microstructure literature have shown the importance of information in determining relations between prices and trade size. For example, Easley and O'Hara (1987) observe that because informed traders prefer to trade larger amounts than uninformed liquidity traders, the size of a trade contains information about who the trader is and, consequently, also contains information about the traders' private information. As a result, prices in their model do not satisfy the Markov property, since the conditional distribution of next period's price depends on the entire history of past prices, i.e., on the order flow. That is, the sequence of price changes of  $1/-1/1$  will have a different effect on the conditional mean than the sequence  $-1/1/1$ , even though both sequences yield the same total price change over the three trades.

One simple implication of such order-flow dependence is that the coefficients of the three lags of  $Z_k$ 's are not identical. If they are, then only the sum of the most recent three price changes matters in determining the conditional mean, and not the order in which those price changes occurred. Therefore, if we denote by  $\beta_p$  the vector of coefficients  $[\beta_2 \ \beta_3 \ \beta_4]'$  of the lagged price changes, the null hypothesis H of order-flow independence is simply:

$$H: \beta_2 = \beta_3 = \beta_4 .$$

This may be recast as a linear hypothesis for  $\beta_p$ , namely  $A\beta_p = 0$ , where

$$A \equiv \begin{pmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \end{pmatrix}. \quad (6.1)$$

Then under  $H_0$ , we obtain the following test statistic:

$$\hat{\beta}'_p A' (A \hat{V}_p A')^{-1} A \hat{\beta}_p \stackrel{a}{\sim} \chi^2_2, \quad (6.2)$$

where  $\hat{V}_p$  is the estimated asymptotic covariance matrix of  $\hat{\beta}_p$ . The values of these test statistics for the six stocks are: IBM = 11,462.43, CUE = 152.05, FWC = 446.01, HNH = 18.62, NAV = 1,184.48, and T = 3,428.92. The null hypothesis of order-flow independence may be rejected at all the usual levels of significance for all six stocks. These findings support Easley and O'Hara's observation that information-based trading can lead to path-dependent price changes, so that the order flow [and the entire history of other variables] may affect the conditional distribution of the next price change.

## 6.2. Measuring price impact per unit volume of trade

By price impact we mean the effect of a current trade of a given size on the conditional distribution of the *subsequent* price change. As such, the coefficients of the variables  $T_\lambda(V_{k-j}) \cdot IBS_{k-j}$ ,  $j = 1, 2, 3$ , measure the price impact of trades per unit of transformed dollar volume. More precisely, recall that our definition of the volume variable is the Box–Cox transformation of dollar volume divided by 100, hence the coefficient  $\beta_{11}$  for stock  $i$  is the contribution to the conditional mean  $X'_k \beta$  that results from a trade of  $\$100 \cdot (1 + \lambda_i)^{1/\lambda_i}$  [since  $T_\lambda((1 + \lambda_i)^{1/\lambda_i}) = 1$ ]. Therefore, the impact of a trade of size  $\$M$  at time  $k - 1$  on  $X'_k \beta$  is simply  $\beta_{11} T_\lambda(M/100)$ . Now the estimated  $\hat{\beta}_{11}$ 's in table 2a are generally positive and significant, with the most recent trade having the largest impact. But this is not the impact we seek, since  $X'_k \beta$  is the conditional mean of the unobserved variable  $Z_k^*$  and not of the observed price change  $Z_k$ . In particular, since  $X'_k \beta$  is scaled by  $\sigma_k$  in (2.10), it is difficult to make meaningful comparisons of the  $\hat{\beta}_{11}$ 's across stocks.

To obtain a measure of a trade's price impact that we *can* compare across stocks, we must translate the impact on  $X'_k \beta$  into an impact on the conditional distribution of the  $Z_k$ 's, conditioned on the trade size and other quantities. Since we have already established that the conditional distribution of price changes is order-flow-dependent, we must condition on a specific sequence of past price

changes and trade sizes. We do this by substituting our parameter estimates into (2.10), choosing particular values for the explanatory variables  $X_k$ , and computing the probabilities explicitly. Specifically, for each stock  $i$  we set  $\Delta t_k$  and  $AB_{k-1}$  to their sample means for that stock and set the remaining regressors to the following values:

$$V_{k-2} = \frac{1}{100} \cdot \text{median dollar volume for stock } i,$$

$$V_{k-3} = \frac{1}{100} \cdot \text{median dollar volume for stock } i,$$

$$SP500_{k-1} = 0.001, \quad SP500_{k-2} = 0.001, \quad SP500_{k-3} = 0.001,$$

$$IBS_{k-1} = 1, \quad IBS_{k-2} = 1, \quad IBS_{k-3} = 1.$$

Specifying values for these variables is equivalent to specifying the market conditions under which price impact is to be measured. These particular values correspond to a scenario in which the most recent three trades are buys, where the sizes of the two earlier trades are equal to the stock's median dollar volume, and where the market has been rising during the past 15 minutes. We then evaluate the probabilities in (2.10) for different values of  $V_{k-1}$ ,  $Z_{k-1}$ ,  $Z_{k-2}$ , and  $Z_{k-3}$ .

For brevity, we focus only on the means of these conditional distributions, which we report for the six stocks in table 3. The entries in the upper panel of table 3 are computed under the assumption that  $Z_{k-1} = Z_{k-2} = Z_{k-3} = +1$ , whereas those in the lower panel are computed under the assumption that  $Z_{k-1} = Z_{k-2} = Z_{k-3} = 0$ . The first entry in the 'IBM' column of table 3's upper panel, -1.315, is the expected price change in ticks of the next transaction of IBM following a \$5,000 buy. The seemingly counterintuitive sign of this conditional mean is the result of the 'bid/ask bounce'; since the past three trades were assumed to be buys, the parameter estimates reflect the empirical fact that the next transaction can be a sell, in which case the transaction price change will often be negative since the price will go from ask to bid. To account for this effect, we would need to include a *contemporaneous* buy/sell indicator,  $IBS_k$ , in  $X'_k$  and condition on this variable as well. But such a variable is clearly endogenous to  $Z_k$  and our parameter estimates would suffer from the familiar simultaneous-equations biases.

In fact, including the contemporaneous buy/sell indicator  $IBS_k$  and contemporaneous transformed volume  $T_\lambda(V_k)$  would yield a more natural measure of price impact, since such a specification, when consistently estimated, can be used to quantify the expected total cost of transacting a given volume. Unfortunately, there are few circumstances in which the contemporaneous buy/sell indicator

Table 3

Price impact of trades as measured by the change in conditional mean of  $Z_k$ , or  $\Delta E[Z_k]$ , when trade sizes are increased incrementally above the base case of a \$5,000 trade. These changes are computed from the ordered probit probabilities for International Business Machines Corporation (IBM - 206,794 trades), Quantum Chemical Corporation (CUE - 26,927 trades), Foster Wheeler Corporation (FWC - 18,199 trades), Handy and Harman Company (HNH - 3,174 trades), Navistar International Corporation (NAV - 96,127 trades), and American Telephone and Telegraph Company (T - 180,726 trades), for the period from 4 January 1988 to 30 December 1988. Price impact measures expressed in percent are percentages of the average of the high and low prices of each security.

	IBM	CUE	FWC	HNH	NAV	T
<b>Increasing price sequence (1/1/1)</b>						
<i>\$Volume</i>						
<b>Price impact in ticks</b>						
E[Z <sub>k</sub> ]: 5,000	- 1.315	- 0.629	- 0.956	- 0.621	- 1.670	- 1.604
ΔE[Z <sub>k</sub> ]: 10,000	0.060	0.072	0.025	0.019	0.017	0.022
ΔE[Z <sub>k</sub> ]: 50,000	0.193	0.239	0.096	0.074	0.070	0.082
ΔE[Z <sub>k</sub> ]: 100,000	0.248	0.310	0.133	0.103	0.100	0.113
ΔE[Z <sub>k</sub> ]: 250,000	0.319	0.403	0.189	0.148	0.148	0.159
ΔE[Z <sub>k</sub> ]: 500,000	0.371	0.473	0.236	0.188	0.191	0.197
<b>Price impact in percent</b>						
E[Z <sub>k</sub> ]: 5,000	- 0.141	- h <sup>b</sup> 0.090	- 0.831	- 0.474	- 3.796	- 0.736
ΔE[Z <sub>k</sub> ]: 10,000	0.006	0.010	0.022	0.015	0.038	0.010
ΔE[Z <sub>k</sub> ]: 50,000	0.021	0.034	0.084	0.057	0.158	0.038
ΔE[Z <sub>k</sub> ]: 100,000	0.027	0.045	0.116	0.079	0.227	0.052
ΔE[Z <sub>k</sub> ]: 250,000	0.034	0.058	0.164	0.113	0.336	0.073
ΔE[Z <sub>k</sub> ]: 500,000	0.040	0.068	0.205	0.143	0.434	0.090
<b>Constant price sequence (0/0/0)</b>						
<i>\$Volume</i>						
<b>Price impact in ticks</b>						
E[Z <sub>k</sub> ]: 5,000	- 0.328	- 0.460	- 0.214	- 0.230	- 0.235	- 0.294
ΔE[Z <sub>k</sub> ]: 10,000	0.037	0.071	0.021	0.018	0.007	0.013
ΔE[Z <sub>k</sub> ]: 50,000	0.120	0.236	0.080	0.070	0.031	0.050
ΔE[Z <sub>k</sub> ]: 100,000	0.155	0.306	0.111	0.098	0.044	0.069
ΔE[Z <sub>k</sub> ]: 250,000	0.200	0.398	0.156	0.140	0.066	0.098
ΔE[Z <sub>k</sub> ]: 500,000	0.234	0.468	0.195	0.177	0.087	0.123
<b>Price impact in percent</b>						
E[Z <sub>k</sub> ]: 5,000	- 0.035	- 0.066	- 0.186	- 0.175	- 0.534	- 0.135
ΔE[Z <sub>k</sub> ]: 10,000	0.004	0.010	0.018	0.014	0.017	0.006
ΔE[Z <sub>k</sub> ]: 50,000	0.013	0.034	0.070	0.053	0.070	0.023
ΔE[Z <sub>k</sub> ]: 100,000	0.017	0.044	0.096	0.074	0.100	0.032
ΔE[Z <sub>k</sub> ]: 250,000	0.021	0.057	0.136	0.107	0.151	0.045
ΔE[Z <sub>k</sub> ]: 500,000	0.025	0.067	0.169	0.135	0.197	0.056

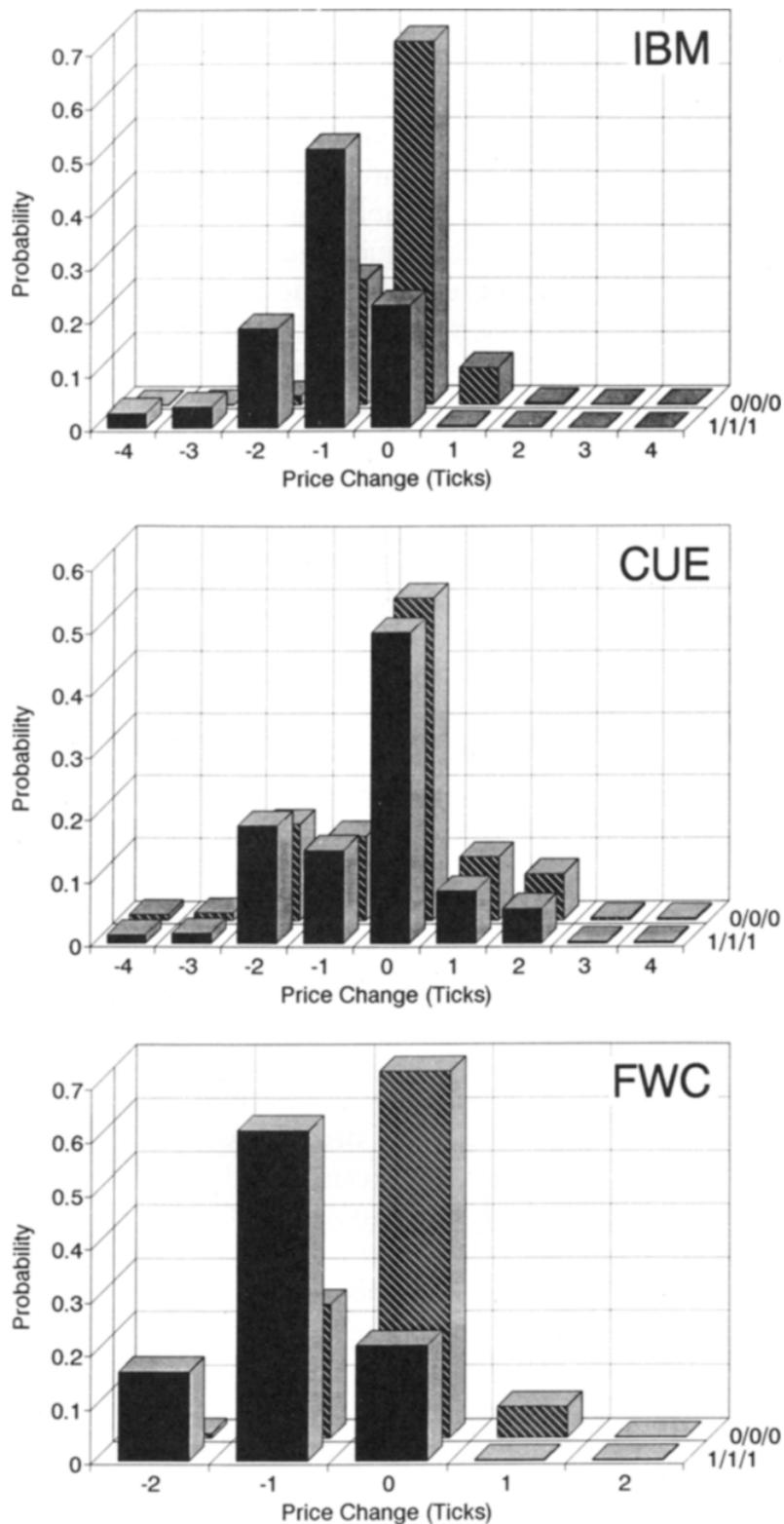
$IBS_k$  may be considered exogenous, since simple economic intuition suggests that factors affecting price changes must also enter the decision to buy or sell. Indeed, limit orders are explicit functions of the current price. Therefore, if  $IBS_k$  is to be included as an explanatory variable in  $X_k$ , its endogeneity must be taken into account. Unfortunately, the standard estimation techniques such as two-stage or three-stage least squares do not apply here because of our discrete dependent variable. Moreover, techniques that allow for discrete dependent variables cannot be applied because the endogenous regressor  $IBS_k$  is also discrete. In principle, it may be possible to derive consistent estimators by considering a joint ordered probit model for both variables, but this is beyond the scope of the current paper. For this reason, we restrict our specification to include only lags of  $IBS_k$  and  $V_k$ .

However, we can 'net out' the effect of the bid/ask spread by computing the change in the conditional mean for trade sizes larger than our base \$5,000 buy. As long as the bid/ask spread remains relatively stable, the change in the conditional mean induced by larger trades will give us a measure of price impact that is independent of it. In particular, the second entry in the 'IBM' column of table 3's upper panel shows that purchasing an additional \$5,000 of IBM (\$10,000 total) increases the conditional mean by 0.060 ticks. However, purchasing an additional \$495,000 of IBM (\$500,000 total) increases the conditional mean by 0.371 ticks; as expected, trading a larger quantity always yields a larger price impact.

A comparison across columns in the upper panel of table 3 shows that larger trades have higher price impact for CUE than for the other five stocks. However, such a comparison ignores the fact that these stocks trade at different price levels, hence a price impact of 0.473 ticks for \$500,000 of CUE may not be as large a percentage of price as a price impact of 0.191 ticks for \$500,000 of NAV. The lower portion of table 3's upper panel reports the price impact as percentages of the average of the high and low prices of each stock, and a trade of \$500,000 does have a higher percentage price impact for NAV than for CUE – 0.434 percent versus 0.068 percent – even though its impact is considerably smaller when measured in ticks. Interestingly, even as a percentage, price impact increases with dollar volume.

In the lower panel of table 3 where price impact values have been computed under the alternative assumption that  $Z_{k-1} = Z_{k-2} = Z_{k-3} = 0$ , the conditional means  $E[Z_k]$  are closer to zero for the \$5,000 buy. For example, the expected price change of NAV is now – 0.235 ticks, whereas in the upper panel it is – 1.670 ticks. Since we are now conditioning on a different scenario, in which the three most recent transactions are buys that have no impact on prices, the empirical estimates imply more probability in the right tail of the conditional distribution of the subsequent price change.

That the conditional mean is still negative may signal the continued importance of the bid/ask spread, nevertheless the price impact measure  $\Delta E[Z_k]$  does



**Fig. 3.** Comparison of estimated ordered probit probabilities of price change, conditioned on a sequence of increasing prices (1/1/1) versus a sequence of constant prices (0/0/0).

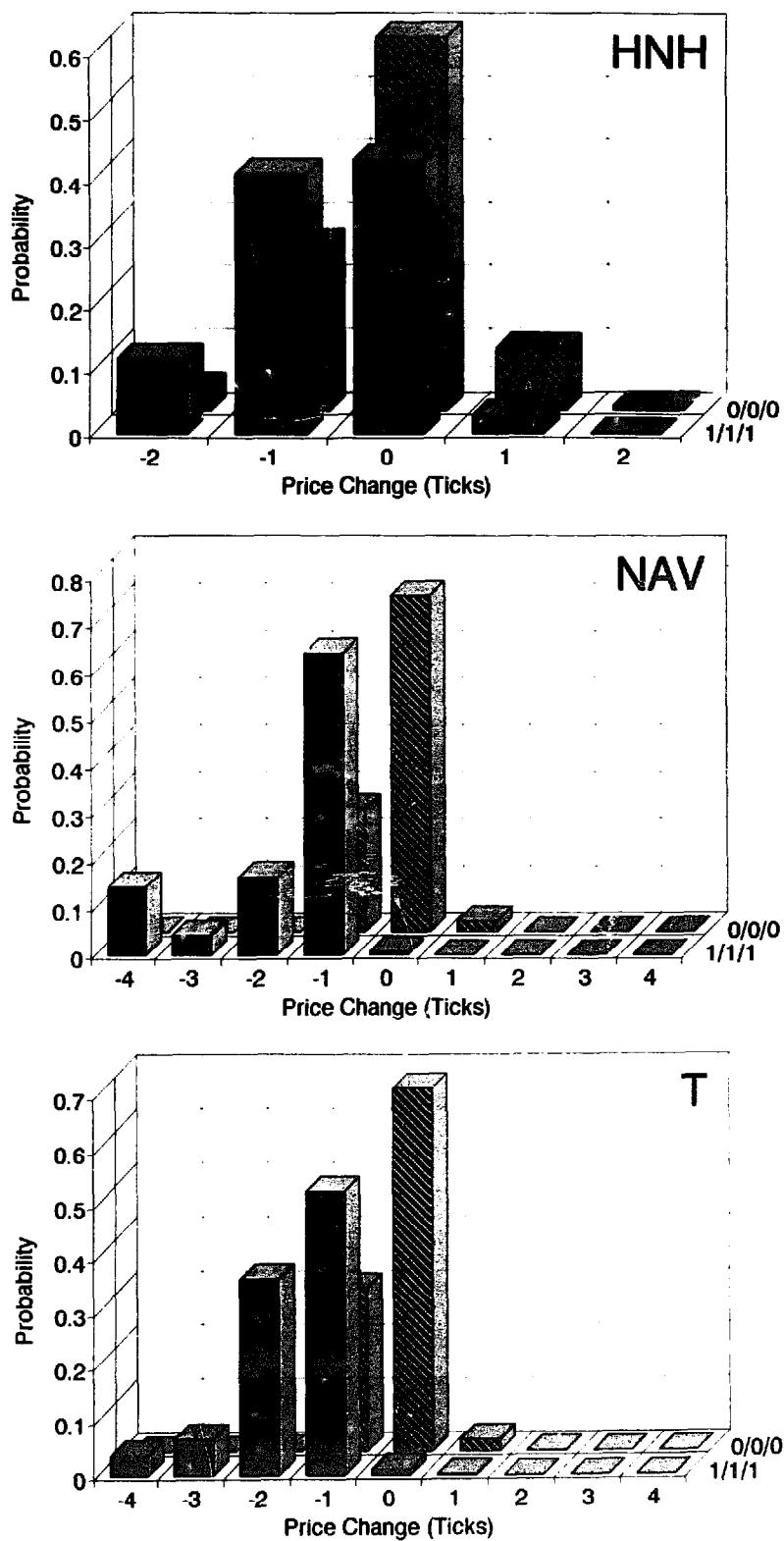


Fig. 3 (continued)

increase with dollar volume in the lower panel. Moreover, these values are similar in magnitude to those in the upper panel – in percentage terms the price impact is virtually the same in both panels of table 3 for most of the six stocks. However, for NAV and T the percentage price impact measures differ considerably between the upper and lower panels of table 3, suggesting that price impact must be measured individually for each security.

Of course, there is no reason to focus solely on the mean of the conditional distribution of  $Z_k$  since we have at our disposal an estimate of the entire distribution. Under the scenarios of the upper and lower panels of table 3 we have also computed the standard deviations of conditional distributions, but since they are quite stable across the two scenarios for the sake of brevity we do not report them here.

To get a sense of their sensitivity to the conditioning variables, we have plotted in fig. 3 the estimated conditional probabilities for the six stocks under both scenarios. In each graph, the cross-hatched bars represent the conditional distribution for the sequence of three buys with a 0 tick price change at each trade, and a fixed trade size equal to the sample median volume for each. The dark-shaded bars represent the conditional distribution for the same sequence of three buys but with a + 1 tick price change for each of the three transactions, also each for a fixed trade size equal to the sample median. The conditional distribution is clearly shifted more to the right under the first scenario than under the second, as the conditional means in table 3 foreshadowed. However, the general shape of the distribution seems rather well-preserved; changing the path of past price changes seems to *translate* the conditional distribution without greatly altering the tail probabilities.

As a final summary of price impact for these securities, we plot ‘price response’ functions in fig. 4 for the six stocks. The price response function, which gives the percentage price impact as a function of dollar volume, reveals several features of the market microstructure that are not as apparent from the numbers in table 3. For example, market liquidity is often defined as the ability to trade any volume with little or no price impact, hence in very liquid markets the price response function should be constant at zero – a flat price response function implies that the percentage price impact is not affected by the size of the trade. Therefore a visual measure of liquidity is the curvature of the price response function; it is no surprise that IBM possesses the flattest price response function of the six stocks.

More generally, the shape of the price response function measures whether there are any economies or diseconomies of scale in trading. An upward-sloping curve implies diseconomies of scale, with larger dollar volume trades yielding a higher percentage price impact. As such, the slope may be one measure of ‘market depth’. For example, if the market for a security is ‘deep’, this is usually taken to mean that large volumes may be traded before much of a price impact is observed. In such cases the price response function may even be

downward-sloping. In fig. 4, all six stocks exhibit trading diseconomies of scale since the price response functions are all upward-sloping, although they increase at a decreasing rate. Such diseconomies of scale suggest that it might pay to break up large trades into sequences of smaller ones. However, recall that the values in fig. 4 are derived from conditional distributions, conditioned on particular sequences of trades and prices. A comparison of the price impact of, say, one \$100,000 trade with two \$50,000 trades can be performed only if the conditional distributions are recomputed to account for the different sequences implicit in the two alternatives. Since these two distinct sequences have not been accounted for in fig. 4, the benefits of dividing large trades into smaller ones cannot be inferred from it. Nevertheless, with the maximum likelihood estimates in hand, such comparisons are trivial to calculate on a case-by-case basis.

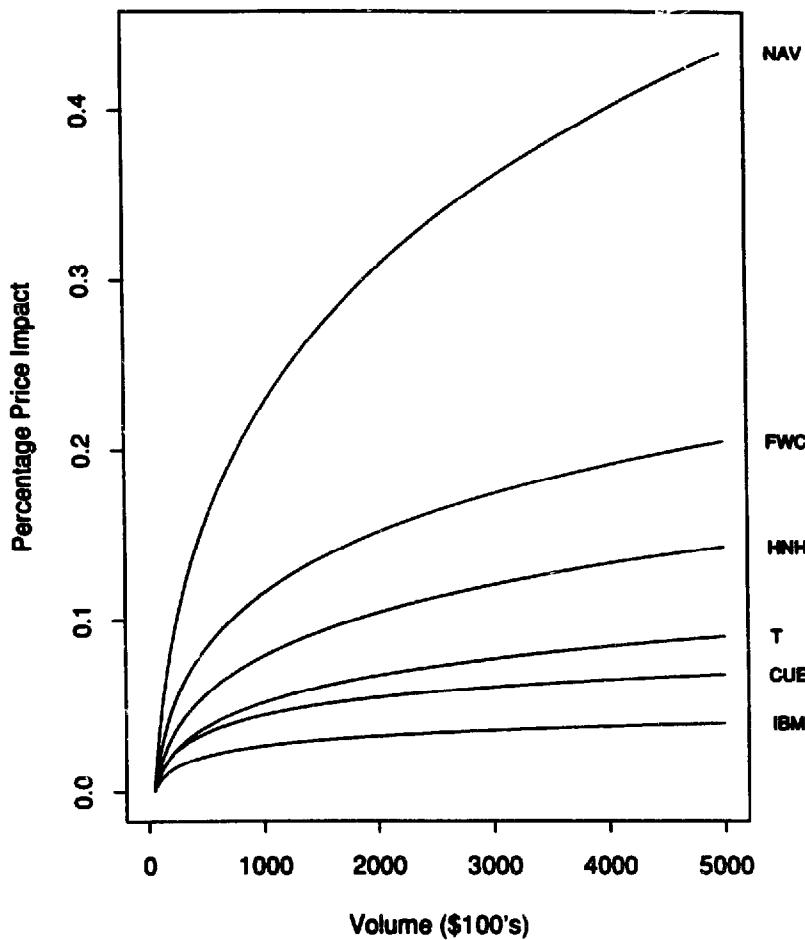


Fig. 4. Percentage price impact as a function of dollar volume computed from ordered probit probabilities, conditional on the three most recent trades being buyer-initiated, and the three most recent price changes being +1 tick each, for IBM (206,794 trades), CUE (26,927 trades), FWC (18,199 trades), HNH (3,174 trades), NAV (96,127 trades), and T (180,726 trades), for the period from 4 January 1988 to 30 December 1988. Percentage price impact is measured as a percentage of the average of the high and low prices for each stock.

Since price response functions are defined in terms of percentage price impact, cross-stock comparisons of liquidity can also be made. Fig. 4 shows that NAV, FWC, and HNH are considerably less liquid than the other stocks, which is partly due to the low price ranges that the three stocks traded in during 1988 (see table 1) – although HNH and NAV have comparable price impacts when measured in ticks (see table 3's upper panel), NAV looks much less liquid when impact is measured as a percentage of price since it traded between \$3.125 and \$7.875, whereas HNH traded between \$14.250 and \$18.500 during 1988. Not surprisingly, since their price ranges are among the highest in the sample, IBM and CUE have the lowest price response functions of the six stocks.

### *6.3. Does discreteness matter?*

Despite the elegance and generality with which the ordered probit framework accounts for price discreteness, irregular trading intervals, and the influence of explanatory variables, the complexity of the estimation procedure raises the question of whether these features can be satisfactorily addressed by a simpler model. Since ordered probit may be viewed as a generalization of the linear regression model to discrete dependent variables, it is not surprising that the latter may share many of the advantages of the former, price discreteness aside. However, linear regression is considerably easier to implement. Therefore, what is gained by ordered probit?

In particular, suppose we ignore the fact that price changes  $Z_k$  are discrete and estimate the following simple regression model via ordinary least squares:

$$Z_k = X'_k \beta + \varepsilon_k. \quad (6.3)$$

Then, suppose we compute the conditional distribution of  $Z_k$  by rounding to the nearest eighth, thus

$$\Pr(Z_k = \frac{j}{8}) = \Pr(\frac{j}{8} - \frac{1}{16} \leq X'_k \beta + \varepsilon_k < \frac{j}{8} + \frac{1}{16}). \quad (6.4)$$

With suitable restrictions on the  $\varepsilon_k$ 's, the regression model (6.3) is known as the 'linear probability' model. The problems associated with applying ordinary least squares to (6.3) are well-known [see for example Judge et al. (1985, ch. 18.2.1)], and numerous extensions have been developed to account for such problems. However, implementing such extensions is at least as involved as maximum likelihood estimation of the ordered probit model and therefore the comparison is of less immediate interest. Despite these problems, we may still ask whether the OLS estimates of (6.3) and (6.4) yield an adequate 'approximation' to a more formal model of price discreteness. Specifically, how different are the probabilities in (6.4) from those of the ordered probit model? If the differences are small, then the linear regression model (6.3) may be an adequate substitute to ordered probit.

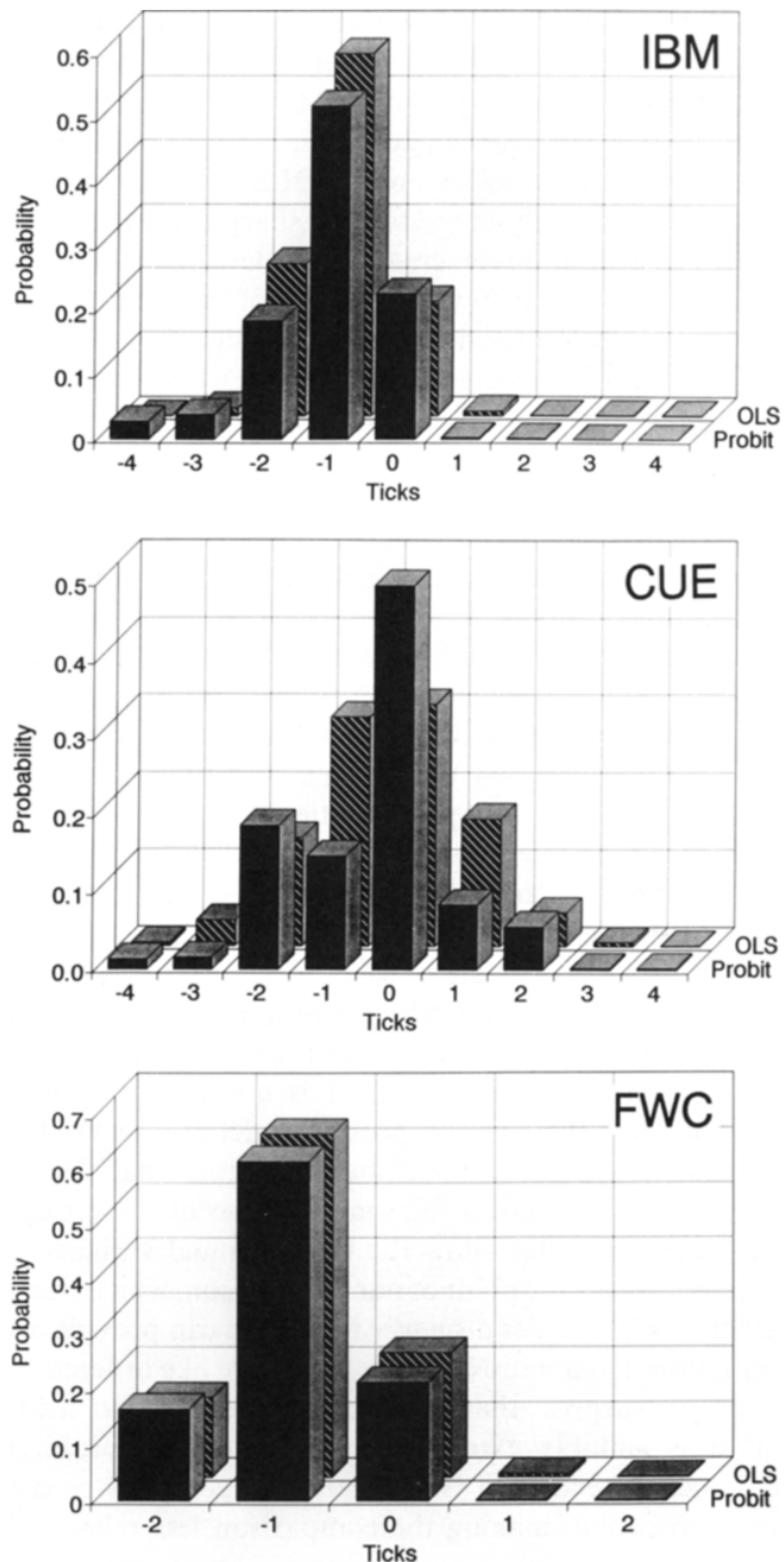
Under the assumption of i.i.d. Gaussian  $\varepsilon_k$ 's, we evaluate the conditional probabilities in (6.4) using the OLS parameter estimates and the same values for the  $X_k$ 's as in section 6.2, and graph them and the corresponding ordered probit probabilities in fig. 5. These graphs show that the two models can yield very different conditional probabilities. All of the OLS conditional distributions are unimodal and have little weight in the tails, in sharp contrast to the much more varied conditional distributions generated by ordered probit. For example, the OLS conditional probabilities show no evidence of the nonmonotonicity that is readily apparent from the ordered probit probabilities of CUE and NAV. In particular, for NAV a price change of -3 ticks is clearly less probable than either -2 or -4 ticks, and for CUE, a price change of -1 tick is less probable than of -2 ticks.

Nevertheless, for FWC the OLS and ordered probit probabilities are rather close. However, it is dangerous to conclude from these matches that OLS is generally acceptable, since these conditional distributions depend sensitively on the values of the conditioning variables. For example, if we plot the same probabilities conditioned on much higher values for  $\sigma_k^2$ , there would be strong differences between the OLS and ordered probit distributions for all six stocks.

Because the ordered probit partition boundaries  $\{\alpha_i\}$  are determined by the data, the tail probabilities of the conditional distribution of price changes may be large or small relative to the probabilities of more central observations, unlike the probabilities implied by (6.3) which are dictated by the (Gaussian) distribution function of . Moreover, it is unlikely that using another distribution function will prov as much flexibility as ordered probit, for the simple reason that (6.3) constrains the state probabilities to be *linear* in the  $X_k$ 's (hence the term 'linear probability model'), whereas ordered probit allows for *nonlinear* effects by letting the data determine the partition boundaries  $\{\alpha_i\}$ .

That OLS and ordered probit can diff. is not surprising given the extra degrees of freedom that the ordered probit model has to fit the conditional distribution of price changes. In fact, it may be argued that the comparison of OLS and ordered probit is not a fair one because of these extra degrees of freedom (for example, why not allow the OLS residual variance to be heteroskedastic?). But this misses the point of our comparison, which was not meant to be fair but rather to see whether a *simpler* technique can provide approximately the same information that a more complex technique like ordered probit does. It should come as no surprise that OLS can come close to fitting nonlinear phenomena if it is suitably extended (in fact, ordered probit is one such extension). But such an extended OLS analysis is generally as complicated to perform as ordered probit, making the comparison less relevant for our purposes.

A more direct test of the difference between ordered probit and the simple 'rounded' linear regression model is to consider the special case of ordered probit in which all the partition boundaries  $\{\alpha_i\}$  are equally spaced and



**Fig. 5. Discreteness matters. A comparison of OLS probabilities versus ordered probit probabilities for price change, conditioned on an increasing price sequence (1/1/1) caused by buyer-initiated trading. Note the nonlinear properties of the CUE and NAV ordered probit probabilities which OLS cannot capture.**

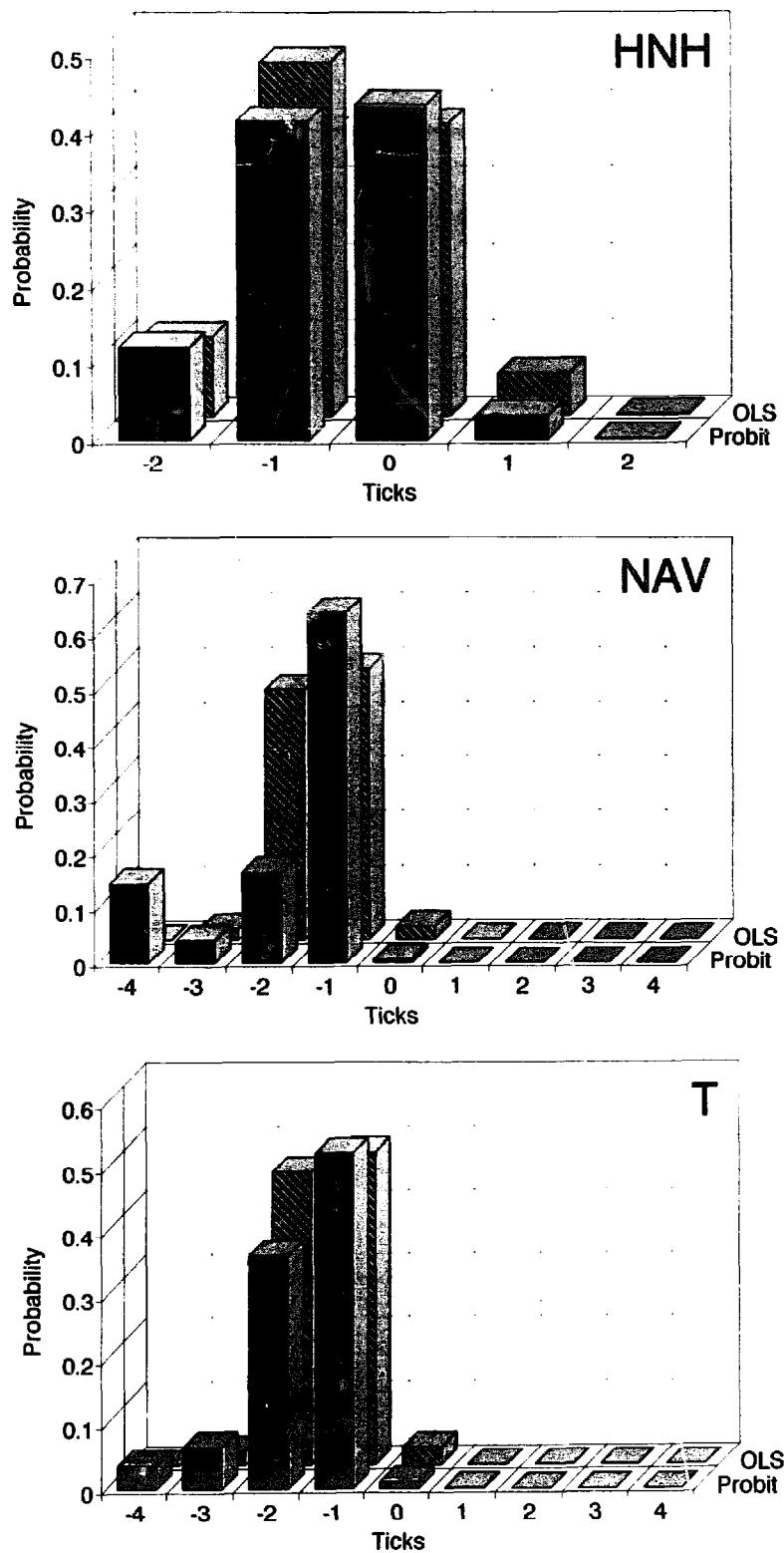


Fig. 5 (continued)

fall on sixteenths. That is, let the observed discrete price change  $Z_k$  be related to the unobserved continuous random variable  $Z_k^*$  in the following manner:

$$Z_k = \begin{cases} -\frac{4}{8} \text{ or less} & \text{if } Z_k^* \in (-\infty, -\frac{4}{8} + \frac{1}{16}), \\ \frac{j}{8} & \text{if } Z_k^* \in [\frac{j}{8} - \frac{1}{16}, \frac{j}{8} + \frac{1}{16}], \quad j = -3, \dots, 3, \\ \frac{4}{8} \text{ or more} & \text{if } Z_k^* \in [\frac{4}{8} - \frac{1}{16}, \infty). \end{cases} \quad (6.5)$$

This is in the spirit of Ball (1988) in which there exists a 'virtual' or 'true' price change  $Z_k^*$  linked to the observed price change  $Z_k$  by rounding  $Z_k^*$  to the nearest multiple of eighths of a dollar. A testable implication of (6.5) is that the partition boundaries  $\{\alpha_i\}$  are equally-spaced, i.e.,

$$\alpha_2 - \alpha_1 = \alpha_3 - \alpha_2 = \dots = \alpha_{m-1} - \alpha_{m-2}, \quad (6.6)$$

where  $m$  is the number of states in our ordered probit model. We can rewrite (6.6) as a linear hypothesis for the  $(m-1) \times 1$ -vector of  $\alpha$ 's in the following way:

$$H: A\alpha = 0, \quad (6.7)$$

where

$$A_{(m-3) \times (m-1)} \equiv \begin{pmatrix} 1 & -2 & 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & -2 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & -2 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 & 1 & -2 & 1 \end{pmatrix}. \quad (6.8)$$

Since the asymptotic distribution of the maximum likelihood estimator  $\hat{\alpha}$  is given by

$$\sqrt{n}(\hat{\alpha} - \alpha) \xrightarrow{a} N(0, \Sigma), \quad (6.9)$$

where  $\Sigma$  is the appropriate submatrix of the inverse of the information matrix corresponding to the likelihood function (2.11), the 'delta method' yields the asymptotic distribution of the following statistic  $\psi$  under the null hypothesis H:

$$H: \psi \equiv n\hat{\alpha}' A'(A\Sigma A')^{-1} A\hat{\alpha} \xrightarrow{a} \chi_{m-3}^2. \quad (6.10)$$

Table 4

Discreteness cannot be completely captured by simple rounding;  $\chi^2$  tests reject the null hypothesis of equally-spaced partition boundaries  $\{\alpha_i\}$  of the ordered probit model for International Business Machines Corporation (IBM – 206,794 trades), Quantum Chemical Corporation (CUE – 26,927 trades), Foster Wheeler Corporation (FWC – 18,199 trades), Handy and Harman Company (HNH – 3,174 trades), Navistar International Corporation (NAV – 96,127 trades), and American Telephone and Telegraph Company (T – 180,726 trades), for the period from 4 January 1988 to 30 December 1988.<sup>a</sup> Entries in the column labeled  $m$  denote the number of states in the ordered probit specification. The 5% and 1% critical values of a  $\chi_2^2$  random variate are 5.99 and 9.21, respectively; the 5% and 1% critical values of a  $\chi_6^2$  random variate are 12.6 and 16.8, respectively.

Stock	Sample size	$\psi \stackrel{a}{\sim} \chi_{m-3}^2$	$m$
IBM	206,794	15,682.35	9
CUE	26,927	366.41	9
FWC	18,199	188.28	5
HNH	3,174	30.59	5
NAV	96,127	998.13	9
T	180,726	1,968.39	9

<sup>a</sup>If price discreteness were simply the result of rounding a continuous 'virtual' price variable to the nearest eighth of a dollar, the ordered probit partition boundaries  $\{\alpha_i\}$  will be equally spaced. If they are, then the statistic  $\psi$  should behave as a  $\chi_{m-3}^2$  variate, where  $m$  is the number of states in the ordered probit specification.

Table 4 reports the  $\psi$ 's for our sample of six stocks, and since the 1% critical values of the  $\chi_2^2$  and  $\chi_6^2$  are 9.21 and 16.8, respectively, we can easily reject the null hypothesis H for each of the stocks. However, because our sample sizes are so large, large  $\chi^2$  statistics need not signal important *economic* departures from the null hypothesis. Nevertheless, the point estimates of the  $\alpha$ 's in table 2a show that they do differ in economically important ways from the simpler rounding model (6.5). With CUE, for example,  $\hat{\alpha}_3 - \hat{\alpha}_2$  is 2.652 but  $\hat{\alpha}_4 - \hat{\alpha}_3$  is 1.031. Such a difference captures the empirical fact that, conditioned on the  $X_k$ 's and  $W_k$ 's, – 1-tick changes are less frequent than – 2-tick changes, even less frequent than predicted by the simple linear probability model.

Discreteness does matter.

## 7. A larger sample

Although our sample of six securities contains several hundred thousand observations, it is still only a small cross-section of the ISSM database, which contains the transactions of over two thousand stocks. It would be impractical for us to estimate our ordered probit model for each one, so we apply our specification to a larger sample of one hundred securities chosen randomly, twenty from each of market-value deciles 6 through 10 (decile 10 contains companies with beginning-of-year market values in the top 10% of the CRSP

Table 5

Names, ticker symbols, market values, and sample sizes over the period from 4 January 1988 to 30 December 1988 for 100 randomly selected stocks for which the ordered probit model was estimated. The selection procedure involved ranking all companies on the CRSP daily returns file by beginning-of-year market value and randomly choosing 20 companies in each of deciles 6 through 10 (decile 10 containing the largest firms), discarding companies which are clearly identified as equity mutual funds. Asterisks next to ticker symbols indicate those securities for which the maximum likelihood estimation procedure did not converge.

Ticker symbol	Company name	Market value x \$1,000	Sample size
<i>Decile 6</i>			
ACP	AMERICAN REAL ESTATE PARTNERS L	217,181	2,394
BCL	BIOCRAFT LABS INC	230,835	7,092
CUL	CULLINET SOFTWARE INC	189,680	18,712
DCY	DCNY CORP	149,073	1,567
FCH	FIRST CAPITAL HLDGS CORP	159,088	8,899
GYK	GIANT YELLOWKNIFE MINES LTD	137,337	1,594
ITX	INTERNATIONAL TECHNOLOGY CORP	161,960	14,675
LOM	LOMAS & NETTLETON MTG INVS	219,450	5,471
MCI*	MASSMUTUAL CORPORATE INVS INC	159,390	727
NET*	NORTH EUROPEAN OIL RTY TR	134,848	708
NPK	NATIONAL PRESTO INDS INC	193,489	1,222
OCQ*	ONEIDA LTD	133,665	1,643
OIL	TRITON ENERGY CORP	195,815	3,203
SII	SMITH INTERNATIONAL INC	148,779	5,435
SKY	SKYLINE CORP	145,821	5,804
SPF	STANDARD PACIFIC CORP DE LP	215,360	11,530
TOL	TOLL BROTHERS INC	157,463	5,519
WIC	WICOR INC	228,044	1,331
WJ	WATKINS JOHNSON CO	192,648	1,647
XTR	XTRA CORP	163,465	1,923
<i>Decile 7</i>			
CER	CILCORP INC	400,138	1,756
CKL	CLARK EQUIPMENT CO	408,509	11,580
CTP	CENTRAL MAINE POWER CO	353,648	5,326
DEI	DIVERSIFIED ENERGIES INC DE	395,505	3,411
FDO	FAMILY DOLLAR STORES INC	286,533	8,513
FRM	FIRST MISSISSIPPI CORP	306,931	8,711
FUR	FIRST UNION REAL EST EQ&MG INVTS	329,041	3,213
KOG	KOGER PROPERTIES INC	265,815	3,508
KWD	KEILWOOD COMPANY	236,271	4,138
LOG	RAYONIER TIMBERLANDS LP	302,500	2,670
MGM	MGMUA COMMUNICATIONS	312,669	10,376
NPR*	NEW PLAN RLTY TR	376,332	1,983
OKE	ONEOK INC	234,668	12,788
SFA	SCIENTIFIC ATLANTA INC	263,801	16,853
SIX*	MOTEL 6 LP	396,768	2,020
SJM	SMUCKER JM CO	373,931	762
SPW	SPX CORP	366,163	7,304
SRR	STRIDE RITE CORP	245,213	5,767
TGR	TIGER INTERNATIONAL INC	352,968	21,612
TRN	TRINITY INDUSTRIES INC	457,366	18,219

Table 5 (continued)

Ticker symbol	Company name	Market value x \$1,000	Sample size
<i>Decile 8</i>			
APS	AMERICAN PRESIDENT COS LTD	617,376	21,554
CAW	CAESARS WORLD INC	525,828	17,900
CBT	CABOT CORP	897,905	5,277
DDS	DILLARD DEPARTMENT STORES INC	758,327	7,267
ERB	ERBAMONT NV	796,698	8,007
FSI	FLIGHT SAFETY INTL INC	833,456	4,562
FVB	FIRST VIRGINIA BANKS INC	496,325	2,637
GLK	GREAT LAKES CHEM CORP	928,358	6,982
HD	HOME DEPOT INC	921,506	16,025
HPH	HARNISCHFEGER INDUSTRIES INC	469,921	7,573
KU	KENTUCKY UTILITIES CO	675,997	8,116
LAC	LAC MINERALS LTD NEW	921,456	4,900
NVP	NEVADA POWER CO	504,785	8,159
ODR	OCEAN DRILLING & EXPL CO	849,965	4,694
PA	PRIMERICA CORP NEW	946,507	35,390
PST	PETRIE STORES CORP	730,688	12,291
REN	ROLLINS ENVIRONMENTAL SVCS INC	825,353	44,272
SW*	STONE & WEBSTER INC	499,568	847
TW	T W SERVICES INC	691,852	16,863
USR	UNITED STATES SHOE CORP	618,686	24,991
<i>Decile 9</i>			
ABS	ALBERTSONS INC	1,695,456	14,171
BDX	BECTON DICKINSON & CO	2,029,188	17,499
CCL	CARNIVAL CRUISE LINES INC	1,294,152	7,111
CYR	CRAY RESEARCH INC	2,180,374	26,459
FFC	FUND AMERICAN COS INC	1,608,525	6,884
FG	USF & G CORP	2,163,821	56,848
GOU	GULF CANADA RESOURCES LIMITED	1,866,365	2,071
GWF	GREAT WESTERN FINANCIAL CORP	1,932,755	20,705
MEA	MEAD CORP	2,131,043	35,796
MEG	MEDIA GENERAL INC	1,002,059	6,304
MLL	MACMILLAN INC	1,387,400	22,083
NSP	NORTHERN STATES POWER CO MN	1,852,777	14,482
PDQ	PRIME MOTOR INNS INC	1,006,803	11,470
PKN	PERKIN ELMER CORP	1,088,400	17,181
RYC	RAYCHEM CORP	1,597,194	16,680
SNG	SOUTHERN NEW ENGLAND TELECOM	1,397,070	4,662
SPS	SOUTHWESTERN PUBLIC SERVICE CO	966,688	10,640
TET	TEXAS EASTERN CORP	1,146,380	29,428
WAG	WALGREEN COMPANY	1,891,310	23,684
WAN	WANG LABS INC	1,801,475	36,607
<i>Decile 10</i>			
AN	AMOCO CORP	7,745,076	39,906
BN	BORDEN INC	3,671,366	22,630
BNI	BURLINGTON NORTHERN INC	4,644,253	33,224
ET	BANKERS TRUST NY CORP	2,426,399	18,502
CAT	CATERPILLAR INC DE	6,137,566	36,379

Table 5 (continued)

Ticker symbol	Company name	Market value x \$1,000	Sample size
CBS	CBS INC	3,709,910	18,630
CCB	CAPITAL CITIES ABC INC	5,581,410	14,585
CPC	CPC INTERNATIONAL INC	3,317,679	27,852
DUK	DUKE POWER CO	4,341,008	17,918
GCI	GANNETT INC	6,335,081	33,512
GIS	GENERAL MILLS INC	4,378,513	26,786
MAS	MASCO CORP	2,867,259	25,746
MHP	MCGRAW HILL INC	2,438,169	36,047
NT	NORTHERN TELECOM LTD	4,049,909	10,128
NYN	NYNEX CORP	3,101,539	40,514
PCG	PACIFIC GAS & ELEC CO	5,982,064	93,981
PFE	PFIZER INC	7,693,452	68,035
RAL	RALSTON PURINA CO	4,517,751	24,710
SGP	SCHEP NG PLOUGH CORP	5,438,652	34,161
UCC	UNION CAMP CORP	2,672,966	14,080

database), also with the restriction that none of these one hundred engaged in stock splits or stock dividends greater than or equal to 3:2. We also discarded (without replacement) randomly chosen stocks that were obviously mutual funds, replacing them with new random draws. Table 5 lists the companies' names, ticker symbols, market values, and number of trades included in our final samples.

Securities from deciles 1 through 5 were not selected because many of them are so thinly traded that the small sample sizes would not permit accurate estimation of the ordered probit parameters. For example, even in deciles 6, 7, and 8, containing companies ranging from \$133 million to \$946 million in market value, there were still six companies for which the maximum likelihood estimation procedure did not converge: MCI, NET, OCQ, NPR, SIX, and SW. In all of these cases, the sample sizes were relatively small, yielding ill-behaved and erratic likelihood functions.

Table 6 presents summary statistics for this sample of one hundred securities broken down by deciles. As expected, the larger stocks tend to have higher prices, shorter times between trades, higher bid/ask spreads (in ticks), and larger median dollar volume per trade. Note that the statistics for  $T_\lambda(V_k) \cdot IBS_k$  implicitly include estimates  $\hat{\lambda}$  of the Box-Cox parameter which differ across stocks. Also, although the mean and standard deviation of  $T_\lambda(V_k) \cdot IBS_k$  for decile 6 differ dramatically from those of the other deciles, these differences are driven solely by the outlier XTR. When this security is dropped from decile 6, the mean and standard deviation of  $T_\lambda(V_k) \cdot IBS_k$  become -0.0244 and 0.3915, respectively, much more in line with the values of the other deciles.

In table 7 we summarize the price impact measures across deciles, where we now define price impact to be the increase in the conditional expected price

change as dollar volume increases from a base case of \$1,000 to either the median dollar volume for each individual stock (the first panel of table 7) or a dollar volume of \$100,000 (the second panel). The first two rows of both panels report decile means and standard deviations of the *absolute* price impact (measured in ticks), whereas the second two rows of both panels report decile means and standard deviations of *percentage* price impact (measured as percentages of the mean of the high and low prices of each stock). For each stock  $i$ , we set  $\Delta t_k$  and  $AB_{k-1}$  to their sample means for that stock and condition on the following values for the other regressors:

$$V_{k-2} = \frac{1}{100} \cdot \text{median dollar volume for stock } i,$$

$$V_{k-3} = \frac{1}{100} \cdot \text{median dollar volume for stock } i,$$

$$SP500_{k-1} = 0.001, \quad SP500_{k-2} = 0.001, \quad SP500_{k-3} = 0.001,$$

$$Z_{k-1} = 1, \quad Z_{k-2} = 1, \quad Z_{k-3} = 1,$$

$$IBS_{k-1} = 1, \quad IBS_{k-2} = 1, \quad IBS_{k-3} = 1,$$

so that we are assuming the three most recent trades are buyer-initiated, accompanied by price increases of one tick each, and the sizes of the two earlier trades are equal to the median dollar volume of the particular stock in question.

From table 7 we see that conditional on a dollar volume equal to the median for the most recent trade, larger capitalization stocks tend to exhibit larger

Table 6

Summary statistics for the sample of 100 randomly chosen securities for the period from 4 January 1988 to 30 December 1988. Market values are computed at the beginning of the year.

Statistic	Deciles				
	6	7	8	9	10
Low price (\$)					
Decile mean	13.94	17.95	21.47	28.02	59.90
Decile std. dev.	9.14	9.75	12.47	12.95	62.27
High price (\$)					
Decile mean	21.11	27.25	33.61	41.39	77.56
Decile std. dev.	11.42	12.16	14.85	21.20	76.93
Market value $\times \$10^9$					
Decile mean	0.177	0.333	0.726	1.602	5.553
Decile std. dev.	0.033	0.065	0.167	0.414	3.737
% prices > midquote					
Decile mean	40.68	41.47	41.77	42.53	43.55
Decile std. dev.	6.36	6.37	3.98	3.71	3.19

Table 6 (continued)

Statistic	Deciles				
	6	7	8	9	10
% prices = midquote					
Decile mean	17.13	19.08	17.91	18.47	16.85
Decile std. dev.	3.99	3.67	4.51	3.93	2.97
% prices < midquote					
Decile mean	42.18	39.45	40.32	39.00	39.60
Decile std. dev.	4.03	4.77	4.30	3.80	2.15
Avg. price change					
Decile mean	0.0085	0.0038	0.0058	-0.0006	0.0015
Decile std. dev.	0.0200	0.0115	0.0103	0.0054	0.0065
Avg. time between trades					
Decile mean	1,085.91	873.66	629.35	430.74	222.49
Decile std. dev.	512.59	489.01	431.79	330.26	109.14
Avg. bid/ask spread					
Decile mean	2.1947	2.3316	2.4926	2.5583	2.9938
Decile std. dev.	0.5396	0.4657	0.3989	0.6514	1.6637
Avg. S&P500 futures return <sup>a</sup>					
Decile mean	-0.0048	-0.0037	-0.0026	-0.0020	-0.0009
Decile std. dev.	0.0080	0.0035	0.0025	0.0019	0.0006
Avg. buy/sell indicator <sup>b</sup>					
Decile mean	-0.0150	0.0202	0.0145	0.0353	0.0395
Decile std. dev.	0.0987	0.1064	0.0695	0.0640	0.0455
Avg. signed transformed volume <sup>c</sup>					
Decile mean	3.9822	0.1969	0.0782	0.2287	0.3017
Decile std. dev.	17.9222	0.6193	0.3230	0.3661	0.2504
Median trading volume (\$)					
Decile mean	6,002	7,345	12,182	16,483	28,310
Decile std. dev.	2,728	3,136	4,985	10,074	13,474
Box-Cox parameter, $\hat{\lambda}^d$					
Decile mean	0.1347	0.0710	0.0127	0.0230	0.0252
Decile std. dev.	0.2579	0.1517	0.0451	0.0679	0.1050

<sup>a</sup>Five-minute continuously-compounded returns of the S&P500 index futures price, for the contract maturing in the closest month beyond the month in which transaction  $k$  occurred, where the return corresponding to the  $k$ th transaction of each stock is computed with the futures price recorded one minute before the nearest round minute *prior* to  $t_k$  and the price recorded five minutes before this.

<sup>b</sup>Takes the value 1 if the  $k$ th transaction price is greater than the average of the quoted bid and ask prices at time  $t_k$ , the value -1 if the  $k$ th transaction price is less than the average of the quoted bid and ask prices at time  $t_k$ , and 0 otherwise.

<sup>c</sup>Box-Cox transformation of dollar volume multiplied by the buy/sell indicator, where the Box-Cox parameter  $\lambda$  is estimated jointly with the other ordered probit parameters via maximum likelihood.

<sup>d</sup>Estimate of Box-Cox parameter  $\lambda$  which determines the degree of curvature that the transformation  $T_\lambda(\cdot)$  exhibits in transforming dollar volume  $V_k$  before inclusion as an explanatory variable in the ordered probit specification. If  $\lambda = 1$ , the transformation  $T_\lambda(\cdot)$  is linear, hence dollar volume enters the ordered probit model linearly. If  $\lambda = 0$ , the transformation is equivalent to  $\log(\cdot)$ , hence the natural logarithm of dollar volume enters the ordered probit model. When  $\lambda$  is between 0 and 1, the curvature of  $T_\lambda(\cdot)$  is between logarithmic and linear.

absolute price impact, no doubt due to their higher prices and their larger median dollar volumes per trade. However, as percentages of the average of their high and low prices, the price impact across deciles is relatively constant as shown by the third row in the first panel of table 7: the average price impact for a median trade in decile 6 is 0.0612%, compared to 0.0523% in decile 10. When conditioning on a dollar volume of \$100,000, however, the results are quite different: the average absolute price impact is similar across deciles, but the average relative price impact is considerably smaller in decile 10 (0.0778%) than in decile 6 (0.2250%). Not surprisingly, a fixed \$100,000 trade will have a greater percentage price impact on smaller capitalization, less liquid stocks than on larger ones.

Further insights on how price impact varies cross-sectionally can be gained from the cross-sectional regressions in table 8, where the four price impact measures and the Box-Cox parameter estimates are each regressed on the following four variables: market value, the initial price level, median dollar volume, and median time-between-trades. Entries in the first row show that the

Table 7

Price impact measures, defined as the increase in conditional expected price change given by the ordered probit model as the volume of the most recent trade is increased from a base case of \$1,000 to either the median level of volume for each security or a level of \$100,000, for the sample of 100 randomly chosen securities for the period from 4 January 1988 to 30 December 1988. Price impact measures expressed in percent are percentages of the average of the high and low prices of each security.

Price impact measure	Deciles				
	6	7	8	9	10
<b>Price impact in ticks</b>					
<b>Lagged volume = Median</b>					
Decile mean	0.0778	0.0991	0.1342	0.1420	0.2020
Decile std. dev.	0.0771	0.0608	0.0358	0.0532	0.0676
<b>Price impact in percent</b>					
<b>Lagged volume = Median</b>					
Decile mean	0.0612	0.0600	0.0703	0.0583	0.0523
Decile std. dev.	0.0336	0.0286	0.0207	0.0229	0.0262
<b>Price impact in ticks</b>					
<b>Lagged volume = \$100,000</b>					
Decile mean	0.2240	0.2611	0.2620	0.2521	0.2849
Decile std. dev.	0.1564	0.1174	0.0499	0.0617	0.0804
<b>Price impact in percent</b>					
<b>Lagged volume = \$100,000</b>					
Decile mean	0.2250	0.1660	0.1442	0.1148	0.0778
Decile std. dev.	0.1602	0.0745	0.0570	0.0633	0.0383

Table 8

Summary of the cross-sectional dispersion in price impact measures and the nonlinearity of the price-change/volume relation (as measured by the Box-Cox parameters,  $\hat{\lambda}_i$ ), via ordinary least-squares regressions for the sample of 100 randomly chosen securities, using market value, initial price, median volume, and median time-between-trades as explanatory variables, for the period from 4 January 1988 to 30 December 1988. Only 94 stocks are included in each of the regressions since the maximum likelihood estimation procedure did not converge for the omitted six. All the coefficients have been multiplied by a factor of 1,000, and z-statistics are given in parentheses, each of which is asymptotically distributed as  $N(0, 1)$  under the null hypothesis that the corresponding coefficient is zero.

Dependent variable	Constant	Market value	Initial price	Median volume	Median $\Delta t_k$	$R^2$
Box-Cox parameter, $\hat{\lambda}_i^a$	118.74 (2.11)	-2.08 (-0.31)	-7.42 (-1.35)	-8.39 (-1.04)	-2.55 (-0.33)	-0.008
Price impact in ticks Lagged volume = Median	93.82 (3.72)	9.86 (3.27)	1.76 (0.71)	5.25 (1.45)	-2.31 (-0.66)	0.184
Price impact in percent Lagged volume = Median	36.07 (4.46)	-1.19 (-1.23)	-2.31 (-2.92)	6.66 (5.72)	0.67 (0.60)	0.376
Price impact in ticks Lagged volume = \$100,000	265.34 (7.03)	8.07 (1.79)	-5.64 (-1.52)	-3.59 (-0.66)	3.25 (0.62)	0.003
Price impact in percent Lagged volume = \$100,000	138.52 (4.17)	-8.53 (-2.15)	-9.61 (-2.95)	8.53 (1.78)	1.74 (0.38)	0.221

<sup>a</sup>The Box-Cox parameter  $\lambda$  determines the degree of curvature that the transformation  $T_\lambda(\cdot)$  exhibits in transforming dollar volume  $V_k$  before inclusion as an explanatory variable in the ordered probit specification. If  $\lambda = 1$ , the transformation  $T_\lambda(\cdot)$  is linear, hence dollar volume enters the ordered probit model linearly. If  $\lambda = 0$ , the transformation is equivalent to  $\log(\cdot)$ , hence the natural logarithm of dollar volume enters the ordered probit model. When  $\lambda$  is between 0 and 1, the curvature of  $T_\lambda(\cdot)$  is between logarithmic and linear.

Box-Cox parameters are inversely related to all four variables, though none of the coefficient estimates are statistically significant and the adjusted  $R^2$  is negative, a symptom of the imprecision with which the  $\lambda_i$ 's are estimated. But the two percentage price impact regressions seem to have higher explanatory power, with adjusted  $R^2$ 's of 37.6% and 22.1%, respectively. These two regressions have identical sign patterns, implying that percentage price impact is larger for smaller stocks, lower-priced stocks, higher-volume stocks, and stocks that trade less frequently.

Of course, these cross-sectional regressions are merely meant as data summaries, and may not correspond to well-specified regression equations. As a further check on the robustness of these regression-based inferences, in table 9 we report Spearman rank correlations between the dependent and independent variables of table 8, which are nonparametric measures of association and are asymptotically normal with mean 0 and variance  $1/(n - 1)$  under the null hypothesis of pairwise independence [see, for example, Randles and Wolfe (1979)]. Since

Table 9

Robust measure of the cross-sectional dispersion in price impact measures and the nonlinearity of the price-change/volume relation (as measured by the Box-Cox parameters,  $\hat{\lambda}_i$ ), via the Spearman rank correlations of  $\hat{\lambda}_i$  and price impact measures with market value, initial price, median volume, and median time-between-trades for the sample of 100 randomly chosen securities, of which 94 are used since the maximum likelihood estimation procedure did not converge for the omitted six, over the period from 4 January 1988 to 30 December 1988. Under the null hypothesis of independence, each of the correlation coefficients is asymptotically normal with mean 0 and variance  $1/(n - 1)$ , hence the two-standard-error confidence interval for these correlation coefficients is  $[-0.207, 0.207]$ .

	Market value	Initial price	Median volume	Median $\Delta t_k$
Box-Cox parameter, $\hat{\lambda}_i^a$	-0.260	-0.503	-0.032	-0.015
Price impact in ticks	0.604	0.678	0.282	-0.360
Lagged volume = Median				
Price impact in percent	-0.156	-0.447	0.486	0.082
Lagged volume = Median				
Price impact in ticks	0.273	0.329	-0.020	-0.089
Lagged volume = \$100,000				
Price impact in percent	-0.547	-0.815	0.088	0.316
Lagged volume = \$100,000				

<sup>a</sup>The Box-Cox parameter  $\lambda$  determines the degree of curvature that the transformation  $T_\lambda(\cdot)$  exhibits in transforming dollar volume  $V_k$  before inclusion as an explanatory variable in the ordered probit specification. If  $\lambda = 1$ , the transformation  $T_\lambda(\cdot)$  is linear, hence dollar volume enters the ordered probit model linearly. If  $\lambda = 0$ , the transformation is equivalent to  $\log(\cdot)$ , hence the natural logarithm of dollar volume enters the ordered probit model. When  $\lambda$  is between 0 and 1, the curvature of  $T_\lambda(\cdot)$  is between logarithmic and linear.

$n = 94$ , the two-standard-error confidence interval about zero for each of the correlation coefficients is  $[-0.207, 0.207]$ . The sign patterns are much the same in table 9 as in table 8, despite the fact that the Spearman rank correlations are not *partial* correlation coefficients.

Such cross-sectional regressions and rank correlations serve only as informal summaries of the data since they are not formally linked to any explicit theories of how price impact should vary across stocks. However, they are consistent with our earlier findings from the six stocks, suggesting that those results are not specific to the behavior of a few possibly peculiar stocks, but may be evidence of a more general and stable mechanism for transaction prices.

## 8. Conclusion

Using 1988 transactions data from the ISSM database, we find that the sequence of trades does affect the conditional distribution for price changes, and the effect is greater for larger capitalization and more actively traded securities. Trade size is also an important factor in the conditional distribution of price changes, with larger trades creating more price pressure, but in a nonlinear

fashion. The price impact of a trade depends critically on the sequence of past price changes and order flows (buy/sell/buy versus sell/buy/buy). The ordered probit framework allows us to compare the price impact of trading over many different market scenarios, such as trading 'with' versus 'against' the market, trading in 'up' and 'down' markets, etc. Finally, we show that discreteness does matter, in the sense that the simpler linear regression analysis of price changes cannot capture all the features of transaction price changes evident in the ordered probit estimates, such as the clustering of price changes on even eighths.

With these applications, we hope to have demonstrated the flexibility and power of the ordered probit model as a tool for investigating the dynamic behavior of transaction prices. Much like the linear regression model for continuous-valued data, the ordered probit model can capture and summarize complex relations between discrete-valued price changes and continuous-valued regressors. Indeed, even in the simple applications considered here, we suffer from an embarrassment of riches in that there are many other empirical implications of our ordered probit estimates that we do not have space to report. For example, we have compared the price impact of only one or two sequences of order flows, price history, and market returns, but there are many other combinations of market conditions, some that might yield considerably different findings. By selecting other scenarios, we may obtain a deeper and broader understanding of how transaction prices react to changing market conditions.

Although we have selected a wide range of regressors to illustrate the flexibility of ordered probit, in practice the specific application will dictate which regressors to include. If, for example, one is interested in testing the implications of Admati and Pfleiderer's (1988) model of intraday patterns in price and volume, natural regressors to include are time-of-day indicators in the conditional mean and variance. If one is interested in measuring how liquidity and price impact vary across markets, an exchange indicator would be appropriate. For intraday event studies, 'event' indicators in both the conditional mean and variance are the natural regressors, and in such cases the generalized residuals we calculated as diagnostics can also be used to construct cumulative average (generalized) residuals.

In the few illustrative applications considered here, we have only hinted at the kinds of insights that ordered probit can yield. The possibilities increase exponentially as we consider the many ways our basic specification can be changed to accommodate the growing number of highly parametrized and less stylized theories about the market microstructure, and we expect to see many other applications in the near future.

## **References**

- Admati, Anat and Paul Pfleiderer, 1988, A theory of intraday patterns: Volume and price variability, *Review of Financial Studies* 1, 3–40.

- Admati, Anat and Paul Pfleiderer, 1989, Divide and conquer: A theory of intraday and day-of-the-week mean effects, *Review of Financial Studies* 2, 189–223.
- Aitchison, James and S. Silvey, 1957, The generalization of probit analysis to the case of multiple responses, *Biometrika* 44, 131–140.
- Amihud, Yakov and Haim Mendelson, 1980, Dealership market: Market-making with inventory, *Journal of Financial Economics* 8, 31–53.
- Amihud, Yakov and Haim Mendelson, 1987, Trading mechanisms and stock returns: An empirical investigation, *Journal of Finance* 42, 533–553.
- Ashford, J., 1959, An approach to the analysis of data for semi-quantal responses in biological response, *Biometrics* 26, 535–581.
- Ball, Clifford, 1988, Estimation bias induced by discrete security prices, *Journal of Finance* 43, 841–865.
- Barclay, Michael and Robert Litzenberger, 1988, Announcement effects of new equity issues and the use of intraday price data, *Journal of Financial Economics* 21, 71–99.
- Berndt, Ernst, Bronwyn Hall, Robert Hall, and Jerry Hausman, 1974, Estimation and inference in nonlinear structural models, *Annals of Economic and Social Measurement* 3, 653–665.
- Blume, Marshall, Craig MacKinlay, and Bruce Terker, 1989, Order imbalances and stock price movements on October 19 and 20, 1987, *Journal of Finance* 44, 827–848.
- Box, George and David Cox, 1964, An analysis of transformations, *Journal of the Royal Statistical Society B* 26, 211–243.
- Bronfman, Corinne, 1991, From trades to orders on the NYSE: Pitfalls in inference using transactions data, Working paper (Department of Finance and Real Estate, College of Business and Public Administration, University of Arizona, Tucson, AZ).
- Campbell, John, Sanford Grossman, and Jiang Wang, 1991, Trading volume and serial correlation in stock returns, Working paper (Princeton University, Princeton, NJ).
- Cho, David and Edward Frees, 1988, Estimating the volatility of discrete stock prices, *Journal of Finance* 43, 451–466.
- Cohen, Kalman, Steven Maier, Robert Schwartz, and David Whitcomb, 1986, The microstructure of securities markets (Prentice-Hall, Englewood Cliffs, NJ).
- Easley, David and Maureen O'Hara, 1987, Price, trade size, and information in securities markets, *Journal of Financial Economics* 19, 69–90.
- Easley, David and Maureen O'Hara, 1990, The process of price adjustment in securities markets, Working paper (Johnson Graduate School of Management, Cornell University, Ithaca, NY).
- Eikeboom, Arnout, 1992, The dynamics of the bid–ask spread, Working paper (Sloan School of Management, Massachusetts Institute of Technology, Cambridge, MA).
- Engle, Robert, 1982, Autoregressive conditional heteroscedasticity with estimates of the variance of United Kingdom inflation, *Econometrica* 50, 987–1007.
- Gallant, Ronald, Peter Rossi, and George Tauchen, 1992, Stock prices and volume, *Review of Financial Studies* 5, 199–242.
- Garman, Mark, 1976, Market microstructure, *Journal of Financial Economics* 3, 257–275.
- Glosten, Lawrence, 1987, Components of the bid–ask spread and the statistical properties of transaction prices, *Journal of Finance* 42, 1293–1307.
- Glosten, Lawrence and Lawrence Harris, 1988, Estimating the components of the bid/ask spread, *Journal of Financial Economics* 21, 123–142.
- Glosten, Lawrence and Paul Milgrom, 1985, Bid, ask and transaction prices in a specialist market with heterogeneously informed traders, *Journal of Financial Economics* 14, 71–100.
- Gottlieb, Gary and Avner Kalay, 1985, Implications of the discreteness of observed stock prices, *Journal of Finance* 40, 135–153.
- Gouriéroux, Christian, Alain Monfort, and Alain Trognon, 1985, A general approach to serial correlation, *Econometric Theory* 1, 315–340.
- Grundy, Bruce and Maureen McNichols, 1989, Trade and the revelation of information through prices and direct disclosure, *Review of Financial Studies* 2, 495–526.
- Gurland, John, Ilbok Lee, and Paul Dahm, 1960, Polychotomous quantal response in biological assay, *Biometrics* 16, 382–398.
- Harris, Lawrence, 1989a, Stock price clustering, discreteness regulation, and bid/ask spreads, Working paper no. 89-01 (New York Stock Exchange, New York, NY).

- Harris, Lawrence, 1989b, Estimation of stock price variances and serial covariances from discrete observations, Working paper (University of Southern California, Los Angeles, CA).
- Harris, Lawrence, 1991, Stock price clustering and discreteness, *Review of Financial Studies* 4, 389–415.
- Harris, Lawrence, George Sofianos, and James Shapiro, 1990, Program trading and intraday volatility, Working paper no. 90-03 (New York Stock Exchange, New York, NY).
- Hasbrouck, Joel, 1988, Trades, quotes, inventories, and information, *Journal of Financial Economics* 22, 229–252.
- Hasbrouck, Joel, 1991a, Measuring the information content of stock trades, *Journal of Finance* 46, 179–207.
- Hasbrouck, Joel, 1991b, The summary informativeness of stock trades: An econometric analysis, *Review of Financial Studies* 4, 571–595.
- Hasbrouck, Joel and Thomas Ho, 1987, Order arrival, quote behavior, and the return-generating process, *Journal of Finance* 42, 1035–1048.
- Hausman, Jerry, 1978, Specification tests in econometrics, *Econometrica* 46, 1251–1271.
- Hausman, Jerry, Andrew Lo, and Craig MacKinlay, 1991, An ordered probit analysis of transaction stock prices, Working paper no. 3888 (NBER, Cambridge, MA).
- Ho, Thomas and Hans Stoll, 1980, On dealer markets under competition, *Journal of Finance* 35, 259–267.
- Ho, Thomas and Hans Stoll, 1981, Optimal dealer pricing under transactions and return uncertainty, *Journal of Financial Economics* 9, 47–73.
- Judge, George, William Griffiths, Carter Hill, Helmut Lütkepohl, and Tsoung-Chao Lee, 1985, *The theory and practice of econometrics* (Wiley, New York, NY).
- Karpoff, Jonathan, 1986, A theory of trading volume, *Journal of Finance* 41, 1069–1087.
- Karpoff, Jonathan, 1987, The relation between price changes and trading volume: A survey, *Journal of Financial and Quantitative Analysis* 22, 109–126.
- Kyle, Albert, 1985, Continuous auctions and insider trading, *Econometrica* 53, 1315–1335.
- Leamer, Edward, 1978, *Specification searches* (Wiley, New York, NY).
- Lee, Charles and Mark Ready, 1991, Inferring trade direction from intraday data, *Journal of Finance* 46, 733–746.
- Lo, Andrew and Craig MacKinlay, 1990a, When are contrarian profits due to stock market overreaction?, *Review of Financial Studies* 3, 175–205.
- Lo, Andrew and Craig MacKinlay, 1990b, Data-snooping biases in tests of financial asset pricing models, *Review of Financial Studies* 3, 431–467.
- Maddala, G., 1983, *Limited-dependent and qualitative variables in econometrics* (Cambridge University Press, Cambridge).
- Madhavan, Ananth and Seymour Smidt, 1991, A Bayesian model of intraday specialist pricing, *Journal of Financial Economics* 30, 99–134.
- McCullagh, Peter, 1980, Regression models for ordinal data, *Journal of the Royal Statistical Society B* 42, 109–127.
- Newey, Whitney, 1985, Semiparametric estimation of limited dependent variable models with endogenous explanatory variables, *Annales de l'INSEE* 59/60, 219–237.
- Petersen, Mitchell, 1986, Testing the efficient market hypothesis: Information lags, the spread, and the role of the market makers, Undergraduate thesis (Princeton University, Princeton, NJ).
- Petersen, Mitchell and Steven Umlauf, 1990, An empirical examination of the intraday behavior of the NYSE specialist, Working paper (Massachusetts Institute of Technology, Cambridge, MA).
- Pritsker, Matthew, 1990, Market microstructure, market efficiency, and the information revealed through trade, Working paper (Princeton University, Princeton, NJ).
- Randles, Ronald and Douglas Wolfe, 1979, *Introduction to the theory of nonparametric statistics* (Wiley, New York, NY).
- Robinson, Michael, 1988, Block trades on the major Canadian and U.S. stock exchanges: A study of pricing behavior and market efficiency, Doctoral dissertation (School of Business Administration, University of Western Ontario, London).
- Roll, Richard, 1984, A simple implicit measure of the effective bid–ask spread in an efficient market, *Journal of Finance* 39, 1127–1139.
- Stoll, Hans, 1989, Inferring the components of the bid–ask spread: Theory and empirical tests, *Journal of Finance* 44, 115–134.

- Stoll, Hans and Robert Whaley, 1990, Stock market structure and volatility, *Review of Financial Studies* 3, 37–71.
- Thisted, Ronald, 1991, Assessing the effect of allergy medications: Models for paired comparisons on ordered categories, *Statistics of Medicine*, forthcoming.
- Wang, Jiang, 1992, A model of competitive stock trading volume, Working paper (Sloan School of Management, Massachusetts Institute of Technology, Cambridge, MA).
- Wood, Robert, Thomas McInish, and Keith Ord, 1985, An investigation of transactions data for NYSE stocks, *Journal of Finance* 40, 723–739.