

# 5

## Hedonic Pricing Models: a Selective and Applied Review

*Stephen Malpezzi*

---

### Introduction

Twenty-five years ago, a few years before his *Housing Economics*, Professor Duncan MacLennan published a seminal paper entitled 'Some Thoughts on the Nature and Purpose of House Price Studies'. The two publications were well-timed: hedonic price modelling was in the process of moving from a cutting-edge empirical curiosity to a standard method of price index construction that has been used in literally thousands of studies since.<sup>1</sup> In those two publications and in MacLennan's subsequent work, several important themes recur:

- the need to put hedonic models on a firm theoretical footing, including, but not limited to, consideration of the consequences of disequilibrium;
- whether common specifications are complete, or reasonably so; including, but not limited to, questions of omitted variables, functional form, and the proper definition of a market; and the implications for our work of inevitably imperfect specifications; and
- that the design of the pricing model should fit the purpose at hand.

In this chapter I selectively review the hedonic price literature (and, briefly, some other related models) with a focus on these questions. Theory will be discussed; but my orientation is more towards the applied economist estimating these models, rather than specialised theorists. In particular, several other recent surveys such as Follain & Jimenez (1985a) and Sheppard (1999) have discussed certain theoretical and econometric issues which I will

discuss only briefly below; readers are referred to those excellent reviews, which I aim to complement rather than compete with.

### **What is a hedonic price index?**

The method of hedonic equations is one way that expenditures on housing can be decomposed into measurable prices and quantities, so that rents for different dwellings or for identical dwellings in different places can be predicted and compared. At its simplest, a hedonic equation is a regression of expenditures (rents or values) on housing characteristics. The independent variables represent the individual characteristics of the dwelling, and the regression coefficients may be transformed into estimates of the implicit prices of these characteristics.

#### *The fundamental hedonic equation*

Hedonic regressions are basically regressions of rent or house value against characteristics of the unit that determine that rent or value. The hedonic regression assumes that one knows the determinants of a unit's rent:

$$R = f(S, N, L, C, T) \quad (5.1)$$

where

R = rent (substitute V, value, if estimating hedonic price indices for, say, homeowners using sales data);

S = structural characteristics;

N = neighbourhood characteristics;

L = location within the market;

C = contract conditions or characteristics, such as whether utilities are included in rent; and

T = the time rent or value is observed.

In this chapter, I will refer to a hedonic model more or less along these lines as a 'single equation' model or the 'first stage' of a 'two-stage' model. Two-stage models attempt to go beyond the initial estimation of a hedonic price

surface, and in the second stage recover structural supply and demand parameters for individual housing characteristics.

Collapsing the vectors  $S$ ,  $N$ ,  $L$  and  $C$  into a larger vector  $X$  for the moment purely for notational convenience, and adopting a common (but sometimes criticised, see below) semi-logarithmic functional form, (5.1) can be re-written compactly as:

$$R = e^{x\beta e} \quad (5.2)$$

so that

$$\ln R = X\beta + \varepsilon \quad (5.3)$$

and we estimate:

$$\ln R = Xb + e \quad (5.4)$$

where  $\beta$  and  $\varepsilon$  are of course the unknown true parameters, and  $b$  and  $e$  are actual estimates.

Now, by properties of logarithms, the predicted rent of a unit can be computed as  $R = e^{xb}$ ; the price of an individual attribute,  $X_1$ , at a given level of  $X_1$ , given the level of the *other*  $m-1$  attributes,  $X_{i \neq 1}$ , can be calculated in dollars or pounds as:

$$P = e^{xb} \quad (5.5)$$

Notice that with such a logarithmic specification, the dollars or pound price of  $X_1$ , or any other single characteristic, varies with the level of  $X_1$ , as well as with the level of other  $X_i$ . Prices are non-linear, an important point to which I return below.

### *'Second-stage' hedonic models: recovering structural parameters*

Much of the hedonic literature focuses on the basic hedonic relationship discussed in the preceding paragraphs. However, as papers discussed below by Rosen and others make clear, the hedonic equation discussed above is a reduced form. Under certain maintained hypotheses hedonic equations also admit of a structural interpretation: for example, if the supply of each and all characteristics is perfectly elastic, hedonic coefficients reveal demand for characteristics. But in most real-world contexts such a stringent maintained

hypothesis is untenable. A number of papers attempt to recover structural parameters of demand and supply; or at least the demand for characteristics.

Specifically, because dollar or pound prices vary within a sample, if the level of characteristics also varies, one can make use of this variation to estimate price elasticities for individual coefficients. For example, one specification that has been used in a number of papers is to presume a linear demand model in the second stage, after a first-stage logarithmic hedonic. It is not uncommon to place prices of each characteristic on the left-hand side, i.e., to assume an inverse demand relation, and to then estimate an equation for each characteristic, of the form:

$$\begin{aligned} P_{1i} &= D_i \alpha_1 + S_i \gamma_1 + \mu_{1i} \\ P_{2i} &= D_i \alpha_2 + S_i \gamma_2 + \mu_{2i} \\ &\vdots \\ P_{mi} &= D_i \alpha_m + S_i \gamma_m + \mu_{mi} \end{aligned} \quad (5.6)$$

where  $P_{1i} = e^{x_i}$  as before. Note that the money price of each of the  $m$  characteristics will vary from one observation to the next (will vary with  $i$ ) because of the property of joint determination of prices discussed above. Armed with this variation in price, demand estimation can proceed.

Papers such as Follain & Jimenez (1985b) and Witte, Sumka & Erekson (1979) present estimates of the demand for housing characteristics from such models. Vectors  $D$  and  $S$  represent exogenous demand and supply shifters, such as income, or input costs (typically land). Often supply is assumed elastic, so that only demand shifters, like household income or family size, that are more readily available, are included. Thus, a prototypical dataset for the first stage would include household level data on some dependent variable like rent or sales price; and on the characteristics  $S$ ,  $N$ ,  $L$  and  $C$ . The prototypical dataset for the second stage would add information on household/unit level demand and supply shifters. Sometimes neighbourhood-level data are appended to household level data.<sup>2</sup>

A point now well understood by most experienced practitioners, but only occasionally discussed in the literature, is just how central a role functional form plays in the set-up for most two-stage demand for characteristics models. Consider that if the first-stage hedonic regression were linear, there would be no variation in characteristic prices within the sample, and hence no second-stage system to estimate. In fact, it is the *difference* between hedonic functional forms, and second-stage demand functional forms, that

makes such systems potentially estimable. This point was made clearly by Nelson's (1982a) insightful critique of Witte, Sumka & Erekson's (1979) otherwise exemplary early study. Witte *et al.* (WSE) estimated logarithmic hedonic *and* logarithmic demand functions; Nelson showed that these were not in fact estimable, and Nelson suggests that the fact that WSE did obtain numerical estimates was, paradoxically, due to rounding error. On the other hand, models that involve (say) logarithmic first-stage hedonic equations and linear second stage demand equations may be estimable (though subject to remaining problems discussed below).

## Repeat sales models

In this chapter, the focus is on hedonic pricing models. In order to place this class of models in context, consider briefly another pricing model, namely repeat sales models, which are in fact related to hedonic price indices in a certain way.<sup>3</sup>

Repeat sales indices are estimated by analysing data where all units have sold at least twice. Such data allow us to annualise the percentage growth in sales prices over time.<sup>4</sup> These are time series indices in their pure form. They do not provide information on the value of individual house characteristics or on price levels. They have the advantage of being based on actual transactions prices, and in principle allow us to sidestep the problem of omitted variable bias.

One way to understand the key features of the repeat sales index is to start by reconsidering the hedonic model. Consider a simple semi-log hedonic equation:

$$\ln P = X\beta + \beta_1 T_1 + \beta_2 T_2 + \beta_3 T_3 + \beta_4 T_4 \quad (5.7)$$

where  $P$  is the value or rent for the unit, and where the vector  $X$  includes all the relevant characters, including a constant term; and the time dummies  $T_i$  represent periods that follow the initial base case period.<sup>5</sup>

The vector  $X$  represents a list of housing and neighbourhood characteristics that would enter a hedonic equation. The vector  $T$  is a series of dummy variables representing the time periods under consideration. These could be months, quarters, or years, depending upon the type of data at hand.

Consider a house 'A' that sells in periods 2 and 4 (period 0 is the base year). In period 2:

$$\begin{aligned}\ln P_2^A &= X\beta + \beta_1 T_1 + \beta_2 T_2 + \beta_3 T_3 + \beta_4 T_4 \\ &= X\beta + \beta_2 T_2\end{aligned}\quad (5.8)$$

since  $T_1$ ,  $T_3$ , and  $T_4 = 0$ . And of course, by similar reasoning, in period 4:

$$\ln P_4^A = X\beta + \beta_4 T_4 \quad (5.9)$$

Then, by subtraction:

$$\begin{aligned}\ln P_4^A - \ln P_2^A &= X\beta + \beta_4 T_4 - X\beta - \beta_2 T_2 \\ &= \beta_4 T_4 - \beta_2 T_2\end{aligned}\quad (5.10)$$

This is for a representative housing unit that sells twice. Given a sample of such units, we want, in effect, the ‘average’  $\beta_4$  and  $\beta_2$ . Recall that regression is, in effect, estimating a series of *conditional means*. Clearly, by subtraction the characteristics vector drops out, as do the dummy variables for periods in which no transaction takes place. Green & Malpezzi (2001) illustrate with sample data.

Another possible refinement is to consider the fact that the variance of these housing prices will generally increase over time. In today’s econometric parlance, such prices are not *stationary*. Case & Shiller (1987) suggest a refinement to the Bailey, Muth and Nourse model to mitigate such problems. The model just described is used as a first stage, and the residuals from this first stage model are used to construct weights that can be used to correct for heteroskedasticity using generalised least squares.

Repeat sales indices are currently much discussed in the literature because they have a number of advantages. First, no information is required on the characteristics of the unit (other than that an individual unit has not significantly changed its characteristics between sales). Second, the method can be used on datasets that are potentially widely available, at least in the US, and collected in a timely manner, with great geographic detail, but which do not have detailed housing characteristics. For example, Case and Shiller’s original work used data collected by the Society of Real Estate Appraisers. Much of the current US research in this area has been undertaken by Fannie Mae and Freddie Mac, which have the advantage of large datasets with price data from a huge number of transactions nationwide.

The repeat sales method has a number of shortcomings as well.<sup>6</sup> First, while raw data have been widely available in the US, data are often harder to come

by in other countries, including Britain. Second, even at its best, the method only yields estimates of price *changes*. No information on price levels, or place-to-place price index, is derivable from the repeat sales method. Of course, the repeat sales method can be combined with some other method, i.e., to update earlier estimates of price levels constructed using some other method. Also, because only a few units transact twice over a given time period, the repeat sales method utilises only a fraction of potential information on the housing market.

Other potential issues with repeat sales include the following. Units that transact frequently may be systematically different from units representative of the stock as a whole. How big a problem this selection bias is depends partly on the purpose of the index. It certainly would be less of a problem if the purpose of the index were to track the prices of units available on the market.

The method also implicitly assumes that there is no change in the quality or quantity of housing services produced by the unit between periods. Of course, this assumption is always violated to some degree. Those who construct these indices spend a lot of time weeding out units that have been upgraded using, for example, collateral data on building permits, or the limited structural information that may exist in the dataset in use. The method also assumes that the coefficients on the underlying hedonic model remain constant: this is what allows the house characteristics to drop out of the model. But this assumption may also be questioned. For example, as families have become smaller, so too has the value of bedrooms, holding all else equal. Thus the hedonic coefficient for bedrooms in 1990 was almost certainly different from the coefficient in 1960, regardless of the particular market.<sup>7</sup>

## The roots of hedonic price models

In essence, the hedonic relation arises because of heterogeneity. The model postulates a market containing a heterogeneous housing stock, which can only be modified at some cost, and heterogeneous consumers, some of whom put different valuations on a given bundle of characteristics ('house') than others. A full history of hedonic pricing models would be another paper, at least, but it is worthwhile to identify the roots of this approach, and point out some of the classics of the literature here.

Two oft-cited classic papers are those by Kelvin Lancaster (1966) and Sherwin Rosen (1974).<sup>8</sup> Focusing on the demand side of the market, Lancaster developed a sophisticated branch of microeconomic theory in which utility

is generated, not by goods *per se*, but by *characteristics* of the goods. The applicability to housing is direct and obvious. I'm happy to be home, not so much to be in anything called a 'house', but to be in a warm dry place, with a quiet space for a comfortable chair, a functioning toilet or a hot bath should I require them, and some other rooms in the house to store stacks of papers or noisy children.<sup>9</sup> Thus, many hedonic studies cite Lancaster's work, and justifiably so, for providing microeconomic foundations for analysing utility-generating characteristics. Lancaster developed this theory using the tools of 'activity analysis', and did not limit his discussion to housing, but applied the concept to topics as diverse as financial assets, the labour/leisure trade-off, and the demand for money. Perhaps Lancaster's main contribution was to put centre stage the still-slippery question of what exactly is the 'good' housing, and how is it related to more fundamental characteristics?

Rosen's (1974) article is the other oft-cited classic reference. Like Lancaster, Rosen focuses on characteristics, but has less to say about their utility-bearing nature and more about how suppliers and consumers interact within a framework of bids and offers for characteristics. Furthermore, while he did not much discuss functional form explicitly, Rosen's model naturally leads to a non-linear hedonic price structure. Many two-stage characteristic demand models, in particular, cite Rosen as their theoretical foundation, although Rosen had little to say about how the estimation of such structural parameters might be carried out.

Of course, there were other early theoretical papers that contributed to the development of this literature. One might point to Sir John Hicks' (1939, 1960) elaboration of a 'composite commodity', and to Fisher & Shell's (1971) related 'repackaging hypothesis', anticipating other early work by Triplett (1974) among others.

As regards the actual estimation of hedonic price models, the study most often cited as the pioneer was not a housing application at all, but a hedonic price index for automobiles developed by A. T. Court (1939). A later, but still early and influential automobile application was by Griliches (1961). Recently, Goodman (1998) and Colwell & Dilmore (1999) have filled in more of this history, telling us more about Court's early work as well as even earlier efforts by Wallace (1926) and Haas (1922). I have seen less discussion of which two-stage study has pride of place; Witte *et al.* (1979) and Awan *et al.* (1982) certainly helped get the ball rolling.



## Conceptual issues in hedonic modelling

In this section, I discuss briefly some matters of underlying theoretical importance in hedonic modelling. The discussion is short because the issues have been so well discussed in two recent surveys, by Follain & Jimenez (1985a), and Sheppard (1999), as well as in several of the other papers cited in the next few paragraphs.

### *Two identification problems*

The identification problem – disentangling supply and demand, when faced with data only from their interaction – has bedevilled applied econometrics for years. In addition to the ‘usual’ identification problem, two-stage hedonic analysis of the demand for characteristics faces an additional potential problem. The problem stems from the nonlinearity of the price structure. In garden-variety supply and demand models, individual consumers (and, often, suppliers) are *price-takers*, i.e., the price of the good is exogenous and the consumer chooses a quantity conditional on the reigning price. In non-linear hedonic models, whether a simple logarithmic model or a more flexible form, prices and quantities are correlated by construction; in effect, consumers choose *both* a quantity of some characteristic and, implicitly, its price. The problem has been well analysed by Blomquist & Worley (1982) and Diamond & Smith (1985), and in the Follain & Jimenez (1985a) and Sheppard (1999) surveys; the reader is referred to their more discursive treatment. Suffice it to say that a number of studies have tried to tackle the problems using well-established techniques like instrumental variables. The problem, as ever, is in finding good instruments.

### *Equilibrium or disequilibrium models?*

Another ubiquitous feature of housing markets is that their extremely costly adjustment processes make the usual simplifying assumptions that markets are in equilibrium when observed less tenable. In fact, the disequilibrium nature of the housing market is a recurring theme in Maclennan (1982) and in many of his other papers. What are the consequences of this disequilibrium, and what can be done about it?

Several possible approaches to issues arising from disequilibrium are possible. In addition to the examples included here, see also the chapter by Yong Tu in this volume. One possible approach to the problem of disequilibrium is to estimate hedonic price functions using only observations in or near

equilibrium. For example, a switching regression approach could be used, following generally the econometric methods described in, for example, Bowden (1978). Such approaches have been applied generally to the housing market, notably by Fair & Jaffee (1972); but their models are generally studies of the determinants of housing starts, not models of housing prices. An example of a later disequilibrium hedonic model is Anas & Eum (1984).

Several difficulties must be overcome to implement such a switching regression model. First, it is necessary to specify the nature of the process that distinguishes equilibrium from disequilibrium observations; in a typical cross-section hedonic sample it is not always obvious how this would be done. Second, depending on the purpose of the hedonic index, successfully estimated equilibrium prices may or may not be what it is needed. For example, if one were constructing place-to-place price indices to study, say, the cost of living or to set appropriate levels of housing subsidies, presumably one would wish the index to reflect actual prices paid in the market, whether in or out of equilibrium.

An alternative approach is to focus, not on the effects of disequilibrium on the construction of the price index, but rather on the amount of disequilibrium in a given sub-market or period, and the adjustment process back to equilibrium. Examples of studies that illustrate the principles include Abraham & Hendershott (1996), Dreiman & Follain (2000), and Malpezzi (1999), each of which studies prices over time. While different in important details, each of these studies essentially proceeds in three steps. First, estimate a time series price index.<sup>10</sup> Second, find some way of differentiating equilibrium prices; in these studies, prices were considered near equilibrium if price *changes* in the next period were near zero. With this subset of equilibrium prices, estimate their fundamental determinants, e.g. based on income, recent growth, supply conditions, etc. Now for each period we have the price and index actually realised, and an estimate of the equilibrium price; thus we have an estimate of disequilibrium as well. The third and final phase of these studies is to study the determinants of the disequilibrium, including in these time-series studies the nature of the time path back to equilibrium, once the market has been 'shocked out'.

## Specification issues

In this section, I discuss a number of practical issues such as which variables or functional form to use, and how to define the geography of a market. But first, one should note an overarching reason why decisions on these matters so often seem ad hoc. It is unfortunate, but the answer to 'what *does* theory

tell us about specification of hedonic models?' is, in brief, 'not much'. Papers like Lancaster (1966) and Rosen (1974) elegantly present models of housing characteristics without having much to say about just what those characteristics are, or how exactly they are related to price.

### *Choice of dependent variable*

Firstly, the choice of dependent variable is about choosing rent, or value, of the housing unit. Confusion sometimes reigns because of sloppy terminology, especially among real estate professionals, but also sometimes among housing economists. The term 'housing price' is often used loosely as a synonym for 'housing value', when of course it is a true 'price' only under special conditions. But the usage is so entrenched by now that even housing economists generally rely on context to keep the meaning clear.

Of course it is well known that house rents and house values are related, though only in special cases proportionately so; papers such as Ambrose and Nourse (1993) and Phillips (1988) explore systematic variations in the relationship, or 'capitalisation rate'. Papers that focus on rents must wrestle with problems stemming from the fact that different units have different lease terms or contract conditions. One notable example is the inclusion, or exclusion, of utility payments in rent. One common procedure is to obtain data on such utility charges for units where they are not included, and to add these charges to contract rent to 'gross it up', so that rents in the sample are for comparable services. Lump sum payments (deposits or 'key money') can be annualised with an assumed capitalisation rate, and added to rent as in Malpezzi's (1998) study of key money paid for Cairo rent controlled units. Another is to use contract rent as the dependent variable, but to add dummy variables to the right-hand side to indicate units with various utilities *included* in rent, so that the estimated coefficients 'dummy out' the price of utilities, leaving a rental index net of utilities.

When estimating a hedonic regression on values, several other measurement issues emerge. A number of studies use owner or tenant estimates of the value of the unit. This gives rise to concern over the accuracy of such self-reported appraisals. Several papers have examined the issue with US data, such as Kain & Quigley (1972b), Follain & Malpezzi (1981b), and Goodman & Ittner (1992). As yet I have found few non-US studies addressing the issue. Early studies such as Kain & Quigley and Follain & Malpezzi suggested that while the *variances* of owner assessments are high, *biases* are modest; given enough data, hedonic models based on owner assessments would be

reasonably reliable. But Goodman and Ittner's recent study finds larger biases, and suggests more caution.

Recent sales 'prices' (house values from observed recent transactions) have some obvious advantages as dependent variables. Recent transactions data may present less potential bias, and greater potential precision, than occupants' or owners' self-assessments. But recent sales are not necessarily a random draw from the total housing stock. If the purpose is to index the market of available units, this may not be of great concern, but if the purpose is to index the total stock, we must concern ourselves with possible selection bias. Several papers such as Gatzlaff & Haurin (1997) have tested the presence of such biases. Test statistics often reject the null, but so far most studies have found the magnitude of the bias to be modest.

Some datasets, like the American Housing Survey, truncate housing values; values over \$300 000 are reported simply as 'over \$300 000'. Such dependent variable truncation can cause significant bias in results. Maddala (1983) presents some econometric techniques that attempt to attenuate the effects of such bias, when truncation is an issue.

### *Selection of independent variables*

There are literally hundreds of potential housing characteristics that could be included on the right-hand side. Butler (1982) and Ozanne and Malpezzi (1985) show that, unfortunately, coefficient estimates are not robust with respect to omitted variables. But interestingly, the same correlation between omitted and included variables that biases individual coefficient estimates can and often does help improve prediction from a 'sparse' model. This suggests that hedonic applications that rely on overall predictions – like place-to-place price indices, or cost-benefit analysis of housing subsidies – can proceed apace, even while papers that rely on interpretation of individual coefficients must be interpreted more cautiously.

While theory is not much of a guide, experience from many studies suggests that, whatever the purpose, a full dataset would include the following:

- rooms, in the aggregate, and by type (bedrooms, bathrooms, etc.)
- floor area of the unit
- structure type (single family, attached or detached, if multi-family the number of units in the structure, number of floors)

- type of heating and cooling systems
- age of the unit
- other structural features, such as presence of basements, fireplaces, garages, etc.
- major categories of structural materials, and quality of finish
- neighbourhood variables, perhaps an overall neighbourhood rating, quality of schools, socio-economic characteristics of the neighbourhood
- distance to the central business district, and perhaps to sub-centres of employment; access to shopping, schools and other important amenities
- characteristics of the tenant that affect prices: length of tenure (especially for renters), whether utilities are included in rent; and possibly racial or ethnic characteristics (if these are hypothesised to affect the price per unit of housing services faced by the occupant)
- date of data collection (especially if the data are collected over a period of months or years).

However, this list, while still incomplete, is also general. Hocking (1976), Amemiya (1980) and Leamer (1978) are among useful guides to the actual selection of variables.

### *Functional form in general*

There is no strong theoretical basis for choosing any specific functional form for a hedonic regression (see Halvorsen & Pollakowski (1981) and Rosen (1974)). Follain & Malpezzi (1980b), for example, tested a linear functional form as well as a log-linear (also known as semi-log) specification. But they found the log-linear form had a number of advantages over the linear form, detailed below.

The log-linear form is written:

$$\ln R = \beta_0 + S\beta_1 + N\beta_2 + L\beta_3 + C\beta_4 + \varepsilon \quad (5.11)$$

where  $\ln R$  is the natural log of imputed rent,  $S$ ,  $N$ ,  $L$  and  $C$  are structural, neighbourhood, locational, and contract characteristics of the dwelling,<sup>11</sup>

and  $\beta_i$  and  $\epsilon$  are the hedonic regression coefficients and error term, respectively.

The log-linear form has five things to recommend it. First, the semi-log model allows for variation in the dollar value of a particular characteristic so that the price of one component depends in part on the house's other characteristics. For example, with the linear model, the value added by a third bathroom to a one-bedroom house is the same as it adds to a five-bedroom house. This seems unlikely<sup>12</sup>. The semi-log model allows the value added to vary proportionally with the size and quality of the home.

Second, the coefficients of a semi-log model have a simple and appealing interpretation. The coefficient can be interpreted as approximately the percentage change in the rent or value given a unit change in the independent variable. For example, if the coefficient of a variable representing central air conditioning is 0.219, then adding it to a structure adds about 22% to its value or its rent. (Actually, the percentage interpretation is an approximation, and it is not necessarily accurate for dummy variables. Halvorsen & Palmquist (1980) show that a much better approximation of the percentage change is given by  $e^b - 1$ , where  $b$  is the estimated coefficient and  $e$  is the base of natural logarithms. So a better approximation is that central air conditioning will add  $\exp(0.219) - 1 = 24\%$ .)

Third, the semi-log form often mitigates the common statistical problem known as heteroskedasticity, or changing variance of the error term. Fourth, semi-log models are computationally simple, and so well suited to examples. The one hazard endemic to the semi-log form is that the anti-log of the predicted log house price does not give an unbiased estimate of predicted price. This can, however, be fixed with an adjustment (see Goldberger 1968). Last, it is possible to build specification flexibility into the right-hand side, using dummy (or indicator) variables, splines and the like (of which more shortly). This allows a fair amount of flexibility in estimation, even with the semi-log form.

However, some authors have recommended more flexible forms than the semi-log. One common flexible form is the translog functional form, suggested by Christensen *et al.* (1973):

$$\ln R = \beta_0 + \sum_m \beta_m \ln X_m + \frac{1}{2} \sum_m \sum_n \gamma_{mn} \ln X_m \ln X_n \quad (5.12)$$

where  $\ln R$  again represents the log of rent (value can be substituted), and there are  $m$  characteristics denoted  $X$ . Examples of the translog form can be found in Capozza *et al.* (1996, 1997).

There is an even more general and flexible class of functions, within which linear, logarithmic and translog functions are subsumed; these flexible forms are carefully analysed by Box and Cox (1964), and applied to hedonic prices by Halvorsen and Pollakowski (1981):

$$R^{\theta} = \beta_0 + \sum_m \beta_m X_m^{\lambda} + \frac{1}{2} \sum_m \sum_n \gamma_{mn} X_m^{\lambda} X_n^{\lambda} \quad (5.13)$$

Such a form is quite flexible, with parameters  $\theta$  and  $\lambda$  limiting the functional form.<sup>13</sup> For example, when  $\theta$  and  $\lambda$  are both 1 and  $\gamma_{mn}$  are all identically zero, the Box-Cox form becomes a simple linear model. When  $\theta$  and  $\lambda$  approach zero and  $\gamma_{mn}$  are all identically zero, the Box-Cox form becomes a logarithmic model. When  $\theta$  and  $\lambda$  approach zero and but some  $\gamma_{mn}$  are nonzero, the Box-Cox form becomes the translogarithmic model.

This is a good place to reiterate the special role functional form plays in two-stage structural models of characteristics demand and supply. I have already noted this important fact: it is functional form – indeed, differences in functional form between stages – that makes the system of demand (or supply and demand) functions potentially estimable. Thus it is particularly problematic that theory yields little guidance to the functional form of the hedonic relationship, and only tenuous guidance to the functional form for second-stage estimation of the demand for characteristics.

### *Functional form and independent variables*

If data permit it, judicious use of dummy or indicator variables for independent variables can be useful. For example, entering a variable for the number of total rooms in (say) a semi-logarithmic hedonic regression constrains the percentage increase in value from a one-unit addition to a 3-room unit to be the same as the percentage increase in value from a one-unit addition to a 6-room unit. If degrees of freedom permit it, at least the most common values can be coded as dummy variables, imparting more flexibility to the form. Malpezzi *et al.* (1980) provide additional details, e.g. how to code combinations of dummies and continuous variables. See also the classic review by Harold Watts (1964), and Suits (1984). The special topic of how to interpret dummy variables when the dependent variable is logarithmic is treated in Halvorsen & Palmquist (1980) and Kennedy (1981).

Of course dummy or indicator variables are not the only method that can be used to incorporate flexibility on the right-hand side. Continuous variables can be entered in quadratic (or cubic or even higher power) form; in fact, as



much flexibility as needed can be readily constructed using piecewise spline techniques (Suits *et al.* 1978).

### *Market and sub-market definition*

The definition and testing of sub-markets is an important recurring theme in MacLennan's work. Housing markets are local and diverse, and hedonic price estimation requires careful consideration of sub-markets (see the chapter in this volume by Yong Tu).

Sub-market assumptions in hedonic models can be roughly categorised as follows. The first category comprises papers that define a market as an entire nation, or at least a large region, or perhaps a state. Linneman (1981) and Struyk (1980) fall into the category of national hedonic models, and Mills & Simenauer (1996) present a regional model. The second category, including much of my own work such as Malpezzi *et al.* (1980) and Follain & Malpezzi (1980a), adopts the metropolitan area as the unit of analysis. Metropolitan areas are usually thought of as labour markets, more or less, and it is certainly appealing to consider housing markets and labour markets as roughly coincident. The third category, including many of MacLennan's own studies as well as papers such as Straszheim (1975), Gabriel (1984), Grigsby *et al.* (1987), Rothenberg *et al.* (1991), MacLennan & Tu (1996), and Bourassa *et al.* (1999), examines sub-markets below the metropolitan level. These may be segmented by location (central city/suburb), or by housing quality level, or by race or income level.

Studies that obtain large datasets and test for the existence of sub-markets, usually by segmenting the sample and performing F-tests for equality of hedonic coefficients across sub-samples, generally find them: the F-tests usually reject the null. Ohta & Griliches (1975) suggest a more conservative method that focuses on changes in the standard error of the regression, in effect on how well the segmented model predicts.

## **Hedonic modelling: the current position**

### *Single-equation models, or structural models?*

Taken together, the problems discussed above – especially the identification problems, imperfect specifications, and the general non-robustness of coefficient estimates – suggest that reliable two-stage structural estimation



of the demand for characteristics is difficult. Qualitatively, that is the judgement we reached in the World Bank's housing demand research project, after investing significant resources attempting to develop characteristic demand models that would improve low-cost housing project design (Follain & Jimenez 1985b; Gross 1986; Mayo & Gross 1987).

That does not mean that there is no hope for developing useful models somewhat along these lines. Models of aggregate housing demand work well enough, despite undoubted problems (see Mayo 1981; Olsen 1987; Whitehead 1999). King (1975) presents an 'in-between' model where housing is broken down into three categories – space, quality and location. Possibly further work along these lines would be fruitful.

### *Hedonic specification: art or science?*

Generally, there is art as well as science in model specification: choice of variables, functional form, and definition of sub-market. Whenever sample sizes are small, and especially if the application will involve some prediction out of sample (as with, say, pricing rent-controlled or subsidised units), it is often best to stick to a simple parsimonious specification, possibly using the metropolitan area as market definition. But if samples are large and well-drawn, and especially if the focus of the hedonic model is a single metropolitan area, more flexible forms, and more careful attention to the delineation of sub-markets, will generally pay off.

It is somewhat surprising that the literature applying formal specification tests, such as those of Hausman (1978), is modest, since specification is such an issue in hedonic analysis. Burgess & Harmon (1982) is an interesting example that could be replicated further.

### **Examples of applications**

While there are many important and interesting theoretical issues related to hedonic models, some of which have been discussed above, our main interest ultimately is in understanding real-world housing markets, and hence in applications. Space precludes an exhaustive review, but here I mention some examples. I list a few representative applications mainly by topic; note also that while I know the US literature best and cite it heavily, hedonic models are now truly universally applied. As previously cited work by Maclennan (1982) and Ball (1973) make clear, there is of course a long-standing, large and growing literature focused on the United Kingdom, as well as the rest of

Europe and of North America; but in fact hedonic models have been applied in every permanently inhabited region of the globe.<sup>14</sup>

One of the first, and still most important, uses of hedonic models is to make general improvements in housing price indices, whether time series, place-to-place, or panel data price indices. Follain & Ozanne (1979), Chowhan & Prud'homme (2000), Englund, Quigley & Redfearn (1998), Follain & Malpezzi (1980b), Hoffman & Kurz (2002), Moulton (1995), Malpezzi, Chun & Green (1998) and Tiwari & Hasegawa (2000) are among many examples of studies that basically aim to improve the precision of housing price benchmarks. Some hedonic studies have been undertaken to construct special-purpose housing price indices, for example to improve the measurement of poverty thresholds (Short *et al.* 1999).

Hedonic prices have also been examined within cities. In addition to the sub-market tests already discussed, many tests of the 'standard urban model' of Alonso, Muth and Mills have been carried out. The standard model predicts a generally declining pattern of prices with distance from the centre of the city. Competing models based on localised amenities, and models with multiple centres, have other predictions. Adair *et al.* (2000), Follain & Malpezzi (1981c), Mozolin (1994), Soderberg & Janssen (2001) are examples of studies that examine intra-urban variation in the price of housing using hedonic models. Perhaps unsurprisingly, results for the 'standard model' are mixed, while there are some broad tendencies for house prices to fall with distance from CBD, amenities and sub-centres generally play an important role as well (also, see the discussion and further references contained in the chapter by Gibb in this volume).

Hedonic models have also been used to develop measures of environmental quality. One common approach is to examine whether house prices increase when near environmental 'goods', or fall when near 'bads'. Examples of this literature include Cheshire & Sheppard (1995), Freeman (1979), Boyle & Kiel (2001), Des Rosiers & Theriault (1996), Din, Hoesli & Bender (2001) and Garrod & Willis (1992 a, b).

Many other interesting studies have been undertaken that focus on interpretations of individual coefficients. One must always be cautious in interpreting individual coefficients in light of the specification issues discussed above. With this *caveat*, a number of studies have examined racial, ethnic and socio-economic differences in housing prices. Kain & Quigley (1972a), Follain & Malpezzi (1981a), Chambers (1992), Galster (1992), Nelson (1982b) and Vandell (1995) are among many contributions to this strand of literature. Other studies have used hedonic age coefficients to measure depreciation,

such as Malpezzi *et al.* (1987), Clapp & Giaccotto (1998b), Goodman & Thibodeau (1995), and Shilling *et al.* (1991).

Hedonic prices have been applied to market-rate units and then used to price subsidised or publicly provided units, in order to calculate the costs and benefits of different housing subsidy programmes. Olsen & Barton (1983), Buchel & Hoesli (1995), De Borger (1986), Quigley (1982a), Satsangi (1991), Turner (1997), Gibb & MacKay (2001) and Willis & Cameron (1993) are representative examples. Closely related are studies that undertake regulatory cost benefit, including rent control; for example, Olsen (1972), Malpezzi (1998) and Willis & Nicholson (1991).

Another important use of hedonic models is the appraisal of individual housing units. Appraisers and other property market professionals increasingly use hedonic models. They can be used to improve professional practice of appraisers and chartered surveyors (Dubin 1998), or for undertaking mass appraisal for property taxation and other public purposes; see Berry & Bednarz (1975), Lusht (1976), and Pace & Gilley (1990).

Hedonic models are also used to examine the capitalisation of a wide range of amenities, as well as costs. One of the earliest literatures along these lines developed to study whether differential local tax rates were capitalised into house prices, following the model of Tiebout (1956), as extended by Oates (1981). After several false starts, papers by Edel & Sclar (1974) and King (1977) clarified the need to include measures of public services as well as taxes paid, and pointed out some important details of the correct functional form for such tests. Many subsequent studies have found such capitalisation, on both the benefit and tax side; Zodrow (1983) provides a convenient review.

Lastly, despite the problems discussed above, many studies have tried to recover demand parameters (and sometimes supply and demand parameters) for individual housing characteristics, or groups of characteristics. Studies here include those by Awan *et al.* (1982), Pasha & Butts (1996), Witte *et al.* (1979) and Kaufman & Quigley (1987).

## Concluding thoughts

A lot of cutting-edge work on conceptual issues relating to hedonics is still being done. Theoretical work continues apace to develop the foundations of hedonic models, and in particular to attempt to address some of the issues I have noted with two-stage structural models. In addition to the literature cited above, on this see, for example, Rouwendal (1992) and Epple (1987).

An alternative way to think about hedonic models is a two-stage process of a different sort: samples used for hedonic estimation are not necessarily random draws from the population of houses, but are selected samples (especially when transactions-driven databases are used). Ermisch *et al.* (1996), Jud & Seaks (1994) and Clapp *et al.* (1991) are examples of studies addressing selection issues.

In terms of functional form, one of the cutting-edge areas is to eschew parametric forms altogether. Semiparametric and nonparametric approaches can be found in Anglin & Ramazan (1996), Mason & Quigley (1996), Meese & Wallace (1991) and Pace (1993). Another approach is to use Bayesian restrictions on hedonic estimates, as outlined in Gilley & Pace (1995) and Knight *et al.* (1992). But perhaps one of the most exciting areas for extending hedonic models is making use of the spatial structure of the data, using the emerging technology of geographic information systems and spatial autocorrelation. Among other recent contributions in this area, see Can (1992), Dubin (1992), Pace & Gilley (1997), Basu & Thibodeau (1998), Gillen *et al.* (2001), and Thibodeau (2002). Thibodeau (2002), for example, finds a roughly 20% improvement in the fit of hedonic models using these techniques. Especially in applications regarding mass appraisal, these techniques are extremely promising.

However, there is also no end of applications that might not be thought of as cutting edge technically, that have not been done, but that are potentially terribly useful. Many of these are extensions of the studies listed in the previous section. Many housing programmes and policies have yet to be submitted to rigorous cost-benefit analysis. Improving systems of mass appraisal remains important; for example, Russia is embarking on the development of a valuation system for all property in the entire country.

While I have already cited individual hedonic studies from every continent, there is clearly scope for more and better international comparisons of housing prices. Recently there has been resurgence in cross-country comparisons, partly driven by the United Nations Centre for Human Settlements' Housing and Urban Development Indicators project.<sup>15</sup> Angel (2000) and Malpezzi & Mayo (1997) present data and comparisons, but these are generally based on simple median house prices from selected cities; more careful analysis, including estimating inter-country hedonic models, remains to be done.<sup>16</sup> An issue of fundamental importance to future hedonic work is the collection of more and better data for hedonic estimation (as well as other kinds of housing analyses). For example, good benchmark data in many countries are still needed. Guides to improved housing data collection include Malpezzi & Mayo (1994) and Malpezzi (2000).

Lastly, while housing is the bulk of every country's real estate, and in fact typically over half a country's tangible capital, hedonic models have rarely been applied to other forms of real estate. Hedonic applications to commercial real estate will be of interest for their own sake; and the functional interdependence of residential and non-residential real estate is often underappreciated by those of us focused on housing.

Over the past three decades, hedonic estimation has clearly matured from a new technology to become the standard way economists deal with housing heterogeneity. Duncan MacLennan's own work in this area, and the work of his colleagues and students, has helped push out these frontiers. Building on this progress, many exciting applications and innovations in hedonic technique undoubtedly lie ahead.

I am grateful to the editors and to fellow contributors to this volume for constructive comments on a previous version. Opinions in this chapter are those of the author, and do not reflect the views of any other individual, or institution.

## Notes

- (1) In drafting this highly selective review, I have benefited from many previous studies, only some of which are listed below. Even the list of other surveys is incomplete. In my own early work, I greatly benefited from Ball's (1973) classic review of early literature. Among many recent reviews, see especially Follain & Jimenez (1985a) and Sheppard (1999).
- (2) Some hedonic studies use aggregate data, e.g. average levels of variables over, say, census tracts; but these have gone out of favour, partly because of aggregation bias (as discussed in Ball 1973) and partly because household/property level data have become more readily available.
- (3) The classic reference on repeat sales is Bailey *et al.* (1963); while there were early applications such as Nourse (1963), they were greatly popularised by several papers by Case & Shiller (e.g. 1987, 1989). Wang & Zorn (1997) provide a thorough review. There are also hybrid models that combine information from both hedonic and repeat sales estimates; see, for example, Case & Quigley (1991) and Quigley (1995). Green & Malpezzi (2001) presents further discussion of these as well as of simpler 'models' such as simple medians of transactions, and Laspeyres, Paasche and Divisia time series indices, and the very important user cost model of price determination.
- (4) Actually, as we will see later in this section, with large samples regression techniques are used, but it amounts to the same thing.

- (5) For notational simplicity we suppress error terms. Careful consideration of house-specific errors and their 'drift' over time is a hallmark of Case & Shiller's (1987) treatment.
- (6) In addition to the references on repeat sales above, see Gatzlaff & Haurin (1997) and Gatzlaff, Green & Ling (1998). I also recommend the excellent review of repeat sales issues contained in Wang & Zorn (1997).
- (7) *Hybrid* indices combine elements of two or more methods into one index. Such methods seek to take advantage of the strengths while minimising the weakness of the constituent indices. These could be time series, cross-section, or both. I have already alluded to hybrid models that combine hedonic and repeat sales methods. The essence of most hybrid models is to 'stack' repeat sales and hedonic models, and then to estimate the two models imposing a constraint that estimated price changes over time are equal in both models. Such methods have the advantage of making use of all available information (see Case & Quigley 1991; Quigley 1995; or Hill *et al.* 1999, for good examples of hybrid indices). Knight *et al.* (1995) use seemingly unrelated regressions as a way to get more efficient coefficient estimates than the coefficient estimates obtained by OLS, but this procedure requires tedious matching of similar observations across years.
- (8) Interestingly, a casual perusal through stacks of papers suggests that UK scholars more often cite Lancaster as their fundamental reference, and US authors more readily cite Rosen. Of course, both papers are often cited by writers from any country. See also Lancaster's later elaboration of his ideas in Lancaster (1971).
- (9) My stepsons view the problem somewhat differently, of course; in their childhood they certainly saw their rooms as refuges from noisy parents.
- (10) Abraham and Hendershott and Follain use repeat sales indices, while Malpezzi combines hedonic and repeat sales indices. But whether hedonic or repeat sales indices are adopted, the general approach is the same.
- (11) Without loss of generality, we've written one of each, when there will usually be several; or if you like, consider each (S, N, L, and C) as a *vector*.
- (12) In fact, if housing units could, via some *Star Trek*-inspired machinery, be instantly and costlessly re-formed, then such costless repackaging would imply a linear structure of prices, where the dollar or pound price of each characteristic was simply added up, much as one would add different items in a shopping basket full of groceries at the checkout to obtain their final price (Fisher & Shell 1971; Triplett 1974). It is ultimately the cost of adjustment that gives rise to the nonlinearities that are observed empirically in housing prices.
- (13) As Halvorsen and Pollakowski point out, additional flexibility could be built in by allowing  $\lambda_m$  to vary with each independent variable; for computational convenience and degrees of freedom, all the hedonic applications Halvorsen and Pollakowski cite, and all those I am familiar with, constrain  $\lambda$  to be the same across all independent variables.
- (14) I have yet to find a hedonic study of Antarctica, but that's about the only region as yet unstudied.

- (15) Initially the indicators project was a joint World Bank–UNCHS research programme, but the World Bank has effectively ceased research in housing and urban development in recent years. See the symposium in *Netherlands Journal of Housing Research* (e.g. Angel, Mayo & Stephens 1993, MacLennan & Gibb, 1993, and Priemus 1992) for critical reviews of the project. See also Flood (1997).
- (16) General discussion of cross-country comparisons include Malpezzi (1990), Annez and Wheaton (1984), Malpezzi and Mayo (1987), MacLennan and Gibb (1993), Harsman and Quigley (1991), Strassman (1991), and Boelhouwer and van der Heijden (1993).