

## **THE "HEDONIC" MODEL USED IN THE ODPM'S HOUSE PRICE INDICES**

### **A description of the regression model developed for the new ODPM House Price Index that was launched in September 2003**

The new monthly House Price Index published by ODPM is based on a multivariate fixed effects regression model ("hedonic" model) of house prices. This paper provides a technical description of the underlying model and how it may be used to estimate average house prices, house price movements and the house price index. For clarity of exposition, the paper contains the following sections:

1. Introduction
  2. Specification of the Hedonic Regression Model
  3. Estimation of the Expected House Price for a Cell
  4. Estimation of the House Price Index
  5. Weighted Regressions
- Appendix - Estimation of Cumulants

## 1. Introduction

The rationale for using the hedonic regression approach is that the number of "cells" (that is, possible combinations of different sets of values for the main effects considered) is so large that it is not possible to estimate accurately, or at all, average house prices for each cell individually from the data available. We assume that the influence of each main effect or interaction is the same across all cells and we include in the model those main effects and interactions which "best" explain variation in house prices. The concept of "best" is determined according to Schwarz's Bayesian Criterion (see Schwarz, 1978 and Shi & Tsai, 2002).

The model used includes the following seven main effects:

- location (local authority district or London borough)
- local authority cluster (an ONS classification of local authorities)
- type of neighbourhood (ACORN)
- dwelling type
- number of habitable rooms (or bedrooms)
- old/new
- first-time buyer/former owner occupier (FTB/FOO)

and the following three interactions:

- ACORN  $\times$  dwelling type
- ACORN  $\times$  FTB/FOO
- dwelling type  $\times$  old/new.

For ease of notation, we denote each cell by a single subscript  $i$  (or  $j$ ). Associated with each cell  $i$  is a vector  $\mathbf{x}_{yi}$  containing the values of the variables in the regression model for cell  $i$ . The subscript  $y$  indicates the calendar year to which the data relate. This is because the new house price index is annually chain-linked and we need to allow for the possibility that the model may change at the January link months (because more housing characteristics become available, for example).

Note that, because all the main effects in the model are categorical, they are represented by dummy, dichotomous covariates and the vector  $\mathbf{x}_{yi}$  contains only the values 1 or 0, according to whether the relevant characteristic is present or not in the cell. Note also that some elements of  $\mathbf{x}_{yi}$  are mutually associated (for example, a dwelling type cannot be both detached and terraced) and all covariates must be linearly independent, to ensure the validity of matrix inversion in estimation. Thus the main effect old/new will be represented by a single covariate with the value (for example) 0 for old and 1 for new.

## 2. Specification of the Hedonic Regression Model

In month  $m$  of year  $y$ , we assume the following model for the price  $P_{ymik}$  of dwelling  $k$  in cell  $i$ :

$$L_{ymik} = \ln(P_{ymik}) = \mathbf{x}'_{yi} \boldsymbol{\beta}_{ym} + \varepsilon_{ymik} \quad 2.1$$

where:  $\boldsymbol{\beta}_{ym}$  is a column vector of  $p_y$  unknown parameters which determine the proportionate impact on the expected price of the  $p_y$  dichotomous covariates included in the column vector  $\mathbf{x}_{yi}$ ;

and the  $\varepsilon_{ymik}$  are independent, identically distributed random variables with the following properties:

$$E[\varepsilon_{ymik}] = 0; \quad \text{Var}[\varepsilon_{ymik}] = \sigma_{ym}^2; \quad \text{Cov}[\varepsilon_{ymik}, \varepsilon_{znjl}] = 0 \quad \text{for } ymik \neq znjl \quad 2.2$$

We also assume that all the higher order central moments of  $\varepsilon_{ymik}$  exist. This is a necessary, but plausible, assumption to cater for the observed circumstance that the  $\varepsilon_{ymik}$  are not Normally distributed. In particular, we assume:

$$E[\varepsilon_{ymik}^3] = \mu_{3ym} = \kappa_{3ym}; \quad E[\varepsilon_{ymik}^4] = \mu_{4ym} = \kappa_{4ym} + 3\sigma_{ym}^4$$

The expected price of a dwelling in cell  $i$ , in month  $m$  of year  $y$  is then:

$$P_{ymi} = E[P_{ymik}] = E[\exp(L_{ymik})] = \exp(\mathbf{x}'_{yi} \boldsymbol{\beta}_{ym}) E[\exp(\varepsilon_{ymik})]$$

From the standard theory of moment generating functions (see, for example, chapter 3 of Stuart & Ord, 1987), we have, for random variable  $x$  and auxiliary variable  $t$ :

$$\ln(E[\exp\{xt\}]) = \sum_{r=1}^{\infty} \frac{\kappa_r t^r}{r!} \quad \text{That is: } E[\exp\{xt\}] = \exp\left(\sum_{r=1}^{\infty} \frac{\kappa_r t^r}{r!}\right) \quad 2.3$$

where  $\kappa_r$  is the  $r^{\text{th}}$  cumulant of the distribution of  $x$ .

For the house price data and models examined,  $\kappa_2 = \sigma^2 \approx 0.15$ ,  $\kappa_3 \approx 0$ ,  $\kappa_4 \approx 0.03$ . Higher order cumulants are negligible and may be ignored.

For  $x = \varepsilon_{ymik}$  and  $t=1$ , we therefore have:

$$E\left[\exp(\varepsilon_{ymik})\right] \approx \exp\left(0 + \frac{\sigma_{ym}^2}{2!} + \frac{\kappa_{3ym}}{3!} + \frac{\kappa_{4ym}}{4!}\right)$$

$$= \exp\left(\frac{\gamma_{ym} \sigma_{ym}^2}{2}\right)$$

$$\text{where } \gamma_{ym} = 1 + \frac{\kappa_{3ym}}{3\sigma_{ym}^2} + \frac{\kappa_{4ym}}{12\sigma_{ym}^2} \quad 2.4$$

is an adjustment factor to correct for non-Normality.

$$\text{Thus: } P_{ymi} = E\left[P_{ymik}\right] = \exp\left(\mathbf{x}'_{yi} \boldsymbol{\beta}_{ym} + \frac{1}{2} \gamma_{ym} \sigma_{ym}^2\right) \quad 2.5$$

### 3. Estimation of the Expected House Price for a Cell

In practice, the vector  $\beta_{ym}$  is not known and has to be estimated. For this purpose, we use the following notation:

$\mathbf{l}_{ym}$  = column vector of all the  $n_{ym}$  elements  $L_{ymik}$  in month  $m$  of year  $y$ ;

$\mathbf{X}_{ym}$  = matrix of housing characteristics for the sample in month  $m$  of year  $y$ :  
each row of the matrix is the row vector  $\mathbf{x}'_{yi}$  of housing characteristics for the corresponding element in vector  $\mathbf{l}_{ym}$ ;

$\boldsymbol{\varepsilon}_{ym}$  = column vector of all the  $n_{ym}$  random components  $\varepsilon_{ymik}$  in month  $m$  of year  $y$ .

Equation 2.1 above relates to a single dwelling. The  $n_{ym}$  equations for month  $m$  of year  $y$  may be summarised in the matrix equation:  $\mathbf{l}_{ym} = \mathbf{X}_{ym}\boldsymbol{\beta}_{ym} + \boldsymbol{\varepsilon}_{ym}$

$$\text{with: } E[\boldsymbol{\varepsilon}_{ym}] = \mathbf{0}; \quad \text{Var}[\boldsymbol{\varepsilon}_{ym}] = \sigma_{ym}^2 \mathbf{I} \quad 3.1$$

Using ordinary least squares estimation, the estimator for  $\boldsymbol{\beta}_{ym}$  is:

$$\hat{\boldsymbol{\beta}}_{ym} = (\mathbf{X}'_{ym}\mathbf{X}_{ym})^{-1} \mathbf{X}'_{ym}\mathbf{l}_{ym} \quad 3.2$$

$$\text{with: } \text{Var}(\hat{\boldsymbol{\beta}}_{ym}) = \sigma_{ym}^2 (\mathbf{X}'_{ym}\mathbf{X}_{ym})^{-1} = \mathbf{V}_{ym} \quad (\text{for notational convenience}) \quad 3.3$$

(see, for example, Johnston, 1972).

In practice, the regression model for the House Price Index uses weighted least squares estimation in order to cater for missing values for some covariates. Section 5 below describes the reasons for this, the treatment applied, the impact on the model and amended formulae for  $\hat{\boldsymbol{\beta}}_{ym}$  and  $\text{Var}(\hat{\boldsymbol{\beta}}_{ym})$ .

We may therefore consider the following estimator for  $P_{ymi}$ :

$$\tilde{P}_{ymi} = \exp\left(\mathbf{x}'_{yi}\hat{\boldsymbol{\beta}}_{ym} + \frac{1}{2}\gamma_{ym}\sigma_{ym}^2\right) \quad 3.4$$

where the unknown vector  $\boldsymbol{\beta}_{ym}$  is replaced by the estimator  $\hat{\boldsymbol{\beta}}_{ym}$  from equation 3.2.

However, because the vector  $\hat{\boldsymbol{\beta}}_{ym}$  is random, the expected value of  $\tilde{P}_{ymi}$  is obtained according to equation 2.3 above. We assume that the scalar random variable  $\mathbf{x}'_{yi}\hat{\boldsymbol{\beta}}_{ym}$  is approximately Normal, because the number of parameters in  $\hat{\boldsymbol{\beta}}_{ym}$  is very much

less than the number of observations on which the estimates are based. With this approximation, only the first two cumulants are non-zero and we have:

$$\begin{aligned} E[\tilde{P}_{ymi}] &= \exp\left(\mathbf{x}'_{yi}\boldsymbol{\beta}_{ym} + \frac{1}{2}\text{Var}\left[\mathbf{x}'_{yi}\hat{\boldsymbol{\beta}}_{ym}\right] + \frac{1}{2}\gamma_{ym}\sigma_{ym}^2\right) \\ &= \exp\left(\mathbf{x}'_{yi}\boldsymbol{\beta}_{ym} + \frac{1}{2}\mathbf{x}'_{yi}\mathbf{V}_{ym}\mathbf{x}_{yi} + \frac{1}{2}\gamma_{ym}\sigma_{ym}^2\right) \quad 3.5 \end{aligned}$$

The estimator  $\tilde{P}_{ymi}$  is therefore biased upwards by the factor  $\exp\left(\frac{1}{2}\mathbf{x}'_{yi}\mathbf{V}_{ym}\mathbf{x}_{yi}\right)$ . An approximately unbiased estimator is therefore:

$$\hat{P}_{ymi} = \exp\left(\mathbf{x}'_{yi}\hat{\boldsymbol{\beta}}_{ym} - \frac{1}{2}\mathbf{x}'_{yi}\mathbf{V}_{ym}\mathbf{x}_{yi} + \frac{1}{2}\gamma_{ym}\sigma_{ym}^2\right) \quad 3.6$$

To estimate  $\sigma_{ym}^2$ ,  $\gamma_{ym}$  and  $\mathbf{V}_{ym}$ , we use sample estimates of  $\sigma_{ym}^2$ ,  $\kappa_{3ym}$  and  $\kappa_{4ym}$ . Because these variance and higher order cumulant terms are also estimated, additional bias correction terms should also, in principle, be included. However, these additional bias correction terms are of the order  $\kappa_{4ym}/n_{ym}$  or less. In practice,  $\kappa_{4ym}$  itself is small and, since  $n_{ym}$  is about 20,000, the additional bias correction term is negligible and may be ignored. Usually,  $\hat{\gamma}_{ym}$  is of the order 1.01.

The appendix to this paper provides formulae, and their derivations, for estimating  $\kappa_{3ym}$  and  $\kappa_{4ym}$ .

#### 4. Estimation of the House Price Index

To calculate the house price index, we calculate for each month a weighted average price for the domain of interest, where the domain may be all dwellings covered by the model or a subset of these dwellings, such as dwellings within a region, old or new dwellings, dwellings of a specific type. We use the estimated mean price in each cell from equation 3.6 and a set of weights  $\{w_{Dyi}\}$ , fixed for each calendar year and

with  $\sum_{i \in D} w_{Dyi} = 1$ , where  $D$  denotes the set of cells included within the domain of interest. The weighted average price for domain  $D$  in month  $m$  of year  $y$  is then:

$$\bar{P}_{Dym} = \sum_{i \in D} w_{Dyi} \hat{P}_{ymi} = \sum_{i \in D} w_{Dyi} \exp\left(\mathbf{x}'_{yi} \hat{\boldsymbol{\beta}}_{ym} - \frac{1}{2} \mathbf{x}'_{yi} \hat{\mathbf{V}}_{ym} \mathbf{x}_{yi} + \frac{1}{2} \hat{\gamma}_{ym} \hat{\sigma}_{ym}^2\right) \quad 4.1$$

where we assume that  $\hat{\mathbf{V}}_{ym}, \hat{\gamma}_{ym}$  and  $\hat{\sigma}_{ym}^2$  are sufficiently good estimators for  $\mathbf{V}_{ym}, \gamma_{ym}$  and  $\sigma_{ym}^2$  that no further bias correction is required (see the discussion after equation 3.6).

The index movement within year  $y$  is  $\bar{P}_{Dym} / \bar{P}_{Dy01}$  and the index relative to February 2002, is:

$$I_{Dym} = 100 \frac{\bar{P}_{D0213}}{\bar{P}_{D0202}} \frac{\bar{P}_{D0313}}{\bar{P}_{D0301}} \dots \frac{\bar{P}_{D,y-1,13}}{\bar{P}_{D,y-1,1}} \frac{\bar{P}_{Dym}}{\bar{P}_{Dy01}} = 100 \left( \frac{\bar{P}_{D0213}}{\bar{P}_{D0202}} \prod_{z=03}^{y-1} \frac{\bar{P}_{Dz13}}{\bar{P}_{Dz01}} \right) \frac{\bar{P}_{Dym}}{\bar{P}_{Dy01}} \quad 4.2$$

where month 13 of year  $y-1$  = month 1 of year  $y$ .

The index may alternatively be expressed as:

$$I_{Dym} = 100 \frac{1}{\bar{P}_{D0202}} \frac{\bar{P}_{D0213}}{\bar{P}_{D0301}} \dots \frac{\bar{P}_{D,y-2,13}}{\bar{P}_{D,y-1,01}} \frac{\bar{P}_{D,y-1,13}}{\bar{P}_{Dy01}} \bar{P}_{Dym} = 100 \left( \prod_{z=02}^{y-1} \frac{\bar{P}_{Dz13}}{\bar{P}_{D,z+1,01}} \right) \frac{\bar{P}_{Dym}}{\bar{P}_{D0202}} \quad 4.3$$

The ratios  $\bar{P}_{Dz13} / \bar{P}_{D,z+1,01}$  are link factors, which ensure continuity of the chain-linked index by correcting for the change, at every January link month, in the weights and model used in the index.

Note that the constraint  $\sum_{i \in D} w_{Dyi} = 1$  is only necessary for the calculation of a weighted average price. Equation 4.2 expresses the house price index as a product of ratios of weighted averages. Each ratio contains the same weights in both numerator and denominator. So, for the house price index, only the relative sizes of the weights are

important and the constraint  $\sum_{i \in D} w_{Dyi} = 1$  is not necessary.

## 5. Weighted Regressions

Sometimes, the data we receive from lenders are incomplete, with missing values for some covariates. Values may be missing for all values of a covariate (or covariates) from a single lender or only for a few observations. Rather than reject observations with missing values, we have chosen to keep them in the regression analysis with a special value for each covariate to indicate a missing value.

In doing so, we are assuming that, for a given set of known covariates, the missing values are "missing at random". That is, the true values for the missing covariates may be represented as a random realisation from the population distribution of these covariates, for the given set of known covariates. This is equivalent to the more general assumption that covariates not included in the model at all (for example: garage; state of repair of the property) are "missing at random" and may be represented by a random error term. Provided this assumption is valid, there is no bias in the parameters estimated in this way.

However, the assumption of constant variance for the error terms is no longer valid because, clearly, observations with missing values will have greater variance than observations with complete covariate data.

To resolve this problem, we assume that the error variance for observations with missing data is equal to the basic error variance  $\sigma_{ym}^2$  for observations with complete data multiplied by a factor greater than 1. The non-uniform variances can then be accommodated by applying a weighted regression, using observation weights equal to the reciprocals of estimates of these factors.

For any given combination of missing covariates, we can estimate these weights, which are less than or equal to 1, from regressions based on complete data only. The estimated observation weight, for a given set of missing covariate values, is calculated as the ratio of the mean squared error from the model including all covariates to the mean squared error from the model including only the non-missing covariates.

It may be desirable to re-calculate these weights every month, using the complete data available for that month. However, this would be time-consuming, when the main operational requirement is to validate the data and produce the current month's index values. We believe that these weights will be stable over time and that small errors in these weights would not have any noticeable impact on the estimated indices. However, these observation weights will be re-assessed annually, to ensure that long-term changes to the factors affecting house prices are picked up and to accommodate any changes to the covariates included in the model.

Using weighted regressions, means that the analyses presented in sections 1 to 3 above need to be modified appropriately. There are no modifications to the underlying principles, because weighted regressions can easily be accommodated by appropriate transformations to the dependent variable and covariates. The required modifications and transformations are covered in standard text books (Johnson, 1972, for example).



In summary, the following equations from sections 1 to 3 need to be modified as shown, to accommodate the application of weighted regressions.

### Sections 1 and 2

The original model in section 1 may be re-written as:

$$L_{ymik} = \ln(P_{ymik}) = \mathbf{x}'_{yi} \boldsymbol{\beta}_{ym} + \varepsilon_{ymik}^* \quad 5.1 \equiv 2.1a$$

where the modified error terms have the properties:

$$E[\varepsilon_{ymik}^*] = 0; \quad \text{Var}[\varepsilon_{ymik}^*] = \frac{\sigma_{ym}^2}{\lambda_{ymik}}; \quad \text{Cov}[\varepsilon_{ymik}^*, \varepsilon_{znjl}^*] = 0 \quad \text{for } ymik \neq znjl \quad 5.2 \equiv 2.2a$$

with  $\lambda_{ymik}$  being the appropriate weight for observation  $ymik$ .

It is possible to transform this model into a model whose error terms have constant variance by transforming the dependent variable and covariates, as follows:

$$\sqrt{\lambda_{ymik}} L_{ymik} = \sqrt{\lambda_{ymik}} \ln(P_{ymik}) = \sqrt{\lambda_{ymik}} \mathbf{x}'_{yi} \boldsymbol{\beta}_{ym} + \varepsilon_{ymik} \quad 5.3 \equiv 2.1b$$

so that the error term  $\varepsilon_{ymik} = \sqrt{\lambda_{ymik}} \varepsilon_{ymik}^*$  has the homoscedastic property of equations 2.2 .

### Section 3

$$E[\boldsymbol{\varepsilon}_{ym}^*] = \mathbf{0}; \quad \text{Var}[\boldsymbol{\varepsilon}_{ym}^*] = \sigma_{ym}^2 \boldsymbol{\Omega}_{ym} = \sigma_{ym}^2 \mathbf{W}_{ym}^{-1} \quad 5.4 \equiv 3.1a$$

where  $\mathbf{W}_{ym}$  is a diagonal matrix of the weights  $\{\lambda_{ymik}\}$ .

Using weighted least squares estimation, the estimator for  $\boldsymbol{\beta}_{ym}$  is:

$$\hat{\boldsymbol{\beta}}_{ym} = (\mathbf{X}'_{ym} \mathbf{W}_{ym} \mathbf{X}_{ym})^{-1} \mathbf{X}'_{ym} \mathbf{W}_{ym} \mathbf{l}_{ym} \quad 5.5 \equiv 3.2a$$

with:  $\text{Var}(\hat{\boldsymbol{\beta}}_{ym}) = \sigma^2 (\mathbf{X}'_{ym} \mathbf{W}_{ym} \mathbf{X}_{ym})^{-1} = \mathbf{V}_{ym}$  (for notational convenience) 5.6  $\equiv$  3.3a

(see, for example, Johnston, 1972, chapter 7).

## References

Cook, J, and Weisberg, S. (1982) *Residuals and Influence in Regression*, Chapman & Hall, London

Johnston, J. (1972) *Econometric Methods*, 2nd ed., McGraw-Hill, Tokyo

Schwarz, G. (1978), Estimating the dimension of a model, *Ann. Statist.* , **6**, 461-464.

Shi, P., Tsai, C-L. (2002), Regression model selection - a residual likelihood approach, *J. R. Statist. Soc. B*, **64**, Part 2, 237-252

Stuart, A. and Ord, J.K. (1987) *Kendall's Advanced Theory of Statistics*, 5th ed., Vol.1, Griffin, London

## Appendix

### Estimation of Cumulants

#### 1. Introduction

The model described in section 1 assumes that the error terms are independent and identically distributed but not that they are Normally distributed. As discussed in section 2, the higher order moments, in the form of cumulants, are relevant for the unbiased estimation of average house prices. This appendix develops estimators for the third and fourth cumulants.

For clarity of exposition, the suffices relating to year, month and cell will be omitted. Suffices will define only the observation number, out of  $n$  observations using a model fitted to  $p$  parameters. For the House Price Index,  $n \approx 20,000$  and  $p \approx 200$ .

The weighted residual from the model for observation  $i$  is defined as:

$$e_i = \sqrt{\lambda_i} \hat{r}_i = \sqrt{\lambda_i} (l_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}) = \sqrt{\lambda_i} l_i - \sqrt{\lambda_i} \mathbf{x}_i' \hat{\boldsymbol{\beta}}$$

Standard regression theory states that the vector of residuals  $\mathbf{e}$  may be expressed as a linear function of the vector of independent, identically distributed error terms  $\boldsymbol{\varepsilon}$  (see, for example, Johnston, 1972):  $\mathbf{e} = (\mathbf{I} - \mathbf{H})\boldsymbol{\varepsilon} = \mathbf{M}\boldsymbol{\varepsilon}$ , where  $\mathbf{I}$  is the identity matrix and  $\mathbf{H}$  is the so-called "hat" matrix  $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$

In particular:

$$e_i = \sum_{j=1}^n M_{ij} \varepsilon_j = \sum_{j=1}^n (I_{ij} - H_{ij}) \varepsilon_j$$

Note that  $I_{jj} = 1$ ,  $I_{ij} = 0$  for  $i \neq j$  and  $(I_{ij})^r = I_{ij}$  for all  $i, j, r$ .

Also:

$$H_{ii} \leq 1, \quad \text{tr}(\mathbf{H}) = \sum_{i=1}^n H_{ii} = p \quad \text{and} \quad \sum_{j=1}^n H_{ij}^2 = H_{ii}$$

because  $\mathbf{H}$  is a symmetric, idempotent matrix of rank  $p$  (see, for example, Cook & Weisberg, 1982).

## 2. Estimating the Error Variance $\sigma^2$

We illustrate our approach to developing estimators for cumulants by applying it to produce the standard unbiased estimator for the error variance  $\sigma^2$ . We have:

$$\begin{aligned}
 E\left[\sum_{i=1}^n e_i^2\right] &= E\left[\sum_{i=1}^n \sum_{j=1}^n M_{ij} \varepsilon_j \sum_{k=1}^n M_{ik} \varepsilon_k\right] \\
 &= \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n M_{ij} M_{ik} E[\varepsilon_j \varepsilon_k] \\
 &= \sum_{i=1}^n \sum_{j=1}^n M_{ij}^2 \sigma^2 \\
 &= \sigma^2 \sum_{i=1}^n \sum_{j=1}^n (I_{ij} - H_{ij})^2 \\
 &= \sigma^2 \sum_{i=1}^n \sum_{j=1}^n (I_{ij} - 2I_{ij}H_{ij} + H_{ij}^2) \\
 &= \sigma^2 \sum_{i=1}^n (1 - 2H_{ii} + H_{ii}) \\
 &= \sigma^2 (n - p)
 \end{aligned}$$

Hence we have the standard, unbiased estimator for the error variance:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{(n - p)}$$

## 3. Estimating the Third Cumulant $\kappa_3$

$$\begin{aligned}
 E\left[\sum_{i=1}^n e_i^3\right] &= E\left[\sum_{i=1}^n \sum_{j=1}^n M_{ij} \varepsilon_j \sum_{k=1}^n M_{ik} \varepsilon_k \sum_{l=1}^n M_{il} \varepsilon_l\right] \\
 &= \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n M_{ij} M_{ik} M_{il} E[\varepsilon_j \varepsilon_k \varepsilon_l] \\
 &= \sum_{i=1}^n \sum_{j=1}^n M_{ij}^3 \mu_3
 \end{aligned}$$

because the independence of the  $\varepsilon_j$  means that  $E[\varepsilon_j \varepsilon_k \varepsilon_l] = 0$  if any subscript  $j, k$  or  $l$  is different from the other two.

Hence, noting also that  $\mu_3 = \kappa_3$ :

$$\begin{aligned}
 E\left[\sum_{i=1}^n e_i^3\right] &= \kappa_3 \sum_{i=1}^n \sum_{j=1}^n (I_{ij} - H_{ij})^3 \\
 &= \kappa_3 \sum_{i=1}^n \sum_{j=1}^n (I_{ij} - 3I_{ij}H_{ij} + 3I_{ij}H_{ij}^2 - H_{ij}^3) \\
 &= \kappa_3 \left\{ \sum_{i=1}^n (1 - 3H_{ii} + 3H_{ii}^2) - \sum_{i=1}^n \sum_{j=1}^n H_{ij}^3 \right\} \\
 &= \kappa_3 \left\{ (n - 3p) + 3 \sum_{i=1}^n H_{ii}^2 - \sum_{i=1}^n \sum_{j=1}^n H_{ij}^3 \right\}
 \end{aligned}$$

The mean value of  $H_{ij}$  is  $p/n$ . If the absolute deviations of the  $H_{ij}$  from the mean  $p/n$  are mainly small (that is, there are not too many large values of  $H_{ij}$ ),

then each  $H_{ij}$  will be of the order  $p/n$  and  $\frac{p^2}{n} \ll \sum_{i=1}^n H_{ii}^2 \leq p$ . In practice, for the

house price data,  $\sum_{i=1}^n H_{ii}^2 \approx 4$ . This is very small relative to  $(n-3p)$  and may be neglected.

Using the same reasoning, the mean value of  $H_{ij}^2$  ( $i \neq j$ ) is approximately

$\frac{p - p^2/n}{n^2 - n} \approx \frac{p}{n^2}$ . Again, if the deviations of  $|H_{ij}|$  from the mean absolute value

$(\sqrt{p})/n$  are sufficiently small, we would expect  $\sum_{i=1}^n \sum_{j=1}^n |H_{ij}^3| \approx \frac{p^{1.5}}{n}$ . In fact, some of

the  $H_{ij}$  may be negative and  $\sum_{i=1}^n \sum_{j=1}^n H_{ij}^3$  will be smaller than this. So this term, too, may be neglected.

We therefore have the following, approximately unbiased estimator for the third cumulant:

$$\hat{\kappa}_3 = \frac{\sum_{i=1}^n e_i^3}{n - 3p}$$

#### 4. Estimating the Fourth Cumulant $\kappa_4$

$$\begin{aligned} E\left[\sum_{i=1}^n e_i^4\right] &= E\left[\sum_{i=1}^n \sum_{j=1}^n M_{ij} \varepsilon_j \sum_{k=1}^n M_{ik} \varepsilon_k \sum_{l=1}^n M_{il} \varepsilon_l \sum_{m=1}^n M_{im} \varepsilon_m\right] \\ &= \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n \sum_{m=1}^n M_{ij} M_{ik} M_{il} M_{im} E[\varepsilon_j \varepsilon_k \varepsilon_l \varepsilon_m] \end{aligned}$$

The independence of the  $\varepsilon_j$  means that  $E[\varepsilon_j \varepsilon_k \varepsilon_l \varepsilon_m] = 0$  if any subscript  $j, k, l$  or  $m$  is different from all the other three. So  $E[\varepsilon_j \varepsilon_k \varepsilon_l \varepsilon_m]$  is not zero only if pairs of subscripts are equal, that is:

$$j=k \text{ and } l=m \ (j \neq l); \ j=l \text{ and } k=m \ (j \neq k); \ j=m \text{ and } k=l \ (j \neq k); \ \text{or } j=k=l=m.$$

The first three conditions are equivalent, corresponding merely to rotation of subscripts, so:

$$\begin{aligned} E\left[\sum_{i=1}^n e_i^4\right] &= \sum_{i=1}^n \sum_{j=1}^n M_{ij}^4 E[\varepsilon_j^4] + 3 \sum_{i=1}^n \sum_{j=1}^n \sum_{\substack{k=1 \\ k \neq j}}^n M_{ij}^2 M_{ik}^2 E[\varepsilon_j^2 \varepsilon_k^2] \\ &= \sum_{i=1}^n \sum_{j=1}^n M_{ij}^4 \mu_4 + 3 \sum_{i=1}^n \sum_{j=1}^n \sum_{\substack{k=1 \\ k \neq j}}^n M_{ij}^2 M_{ik}^2 E[\varepsilon_j^2] E[\varepsilon_k^2] \quad (\text{from the independence of } \varepsilon_j \text{ and } \varepsilon_k) \\ &= \sum_{i=1}^n \sum_{j=1}^n M_{ij}^4 (\kappa_4 + 3\sigma^4) + 3 \sum_{i=1}^n \sum_{j=1}^n \sum_{\substack{k=1 \\ k \neq j}}^n M_{ij}^2 M_{ik}^2 \sigma^4 \\ &= \sum_{i=1}^n \sum_{j=1}^n M_{ij}^4 \kappa_4 + 3 \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n M_{ij}^2 M_{ik}^2 \sigma^4 \\ &= \sum_{i=1}^n \sum_{j=1}^n M_{ij}^4 \kappa_4 + 3 \sum_{i=1}^n \left( \sum_{j=1}^n M_{ij}^2 \right)^2 \sigma^4 \end{aligned}$$

Hence:

$$\begin{aligned}
E\left[\sum_{i=1}^n e_i^4\right] &= \kappa_4 \sum_{i=1}^n \sum_{j=1}^n (I_{ij} - H_{ij})^4 + 3\sigma^4 \sum_{i=1}^n \left\{ \sum_{j=1}^n (I_{ij} - H_{ij})^2 \right\}^2 \\
&= \kappa_4 \sum_{i=1}^n \sum_{j=1}^n (I_{ij} - 4I_{ij}H_{ij} + 6I_{ij}H_{ij}^2 - 4I_{ij}H_{ij}^3 + H_{ij}^4) + 3\sigma^4 \sum_{i=1}^n \left\{ \sum_{j=1}^n (I_{ij} - 2I_{ij}H_{ij} + H_{ij}^2) \right\}^2 \\
&= \kappa_4 \left\{ \sum_{i=1}^n (1 - 4H_{ii} + 6H_{ii}^2 - 4H_{ii}^3) + \sum_{i=1}^n \sum_{j=1}^n H_{ij}^4 \right\} + 3\sigma^4 \sum_{i=1}^n \{1 - H_{ii}\}^2 \\
&= \kappa_4 \left\{ n - 4p + 6 \sum_{i=1}^n H_{ii}^2 - 4 \sum_{i=1}^n H_{ii}^3 + \sum_{i=1}^n \sum_{j=1}^n H_{ij}^4 \right\} + 3\sigma^4 \sum_{i=1}^n (1 - 2H_{ii} + H_{ii}^2) \\
&= \kappa_4 \left\{ n - 4p + 6 \sum_{i=1}^n H_{ii}^2 - 4 \sum_{i=1}^n H_{ii}^3 + \sum_{i=1}^n \sum_{j=1}^n H_{ij}^4 \right\} + 3\sigma^4 \left( n - 2p + \sum_{i=1}^n H_{ii}^2 \right)
\end{aligned}$$

Using the same arguments presented above for the third cumulant  $\kappa_3$ , we may ignore the summation terms in the line above to obtain the approximation:

$$E\left[\sum_{i=1}^n e_i^4\right] \approx \kappa_4 \{n - 4p\} + 3\sigma^4 (n - 2p)$$

from which we have the following, approximately unbiased estimator for the fourth cumulant  $\kappa_4$ :

$$\hat{\kappa}_4 = \frac{\sum_{i=1}^n e_i^4 - 3\hat{\sigma}^4 (n - 2p)}{(n - 4p)}$$

Note that  $\hat{\sigma}^4 = (\hat{\sigma}^2)^2$  is a biased estimator of  $\sigma^4$  but the bias is of the order  $\kappa_4/n$ , which is negligible to the degree of approximation used here.