

**CENTER ON JAPANESE ECONOMY AND BUSINESS**

---

日本経済経営研究所

Working Paper Series

May 2011, No. 296

---

# **On the Evolution of the House Price Distribution**

Takaaki Ohnishi, Takayuki Mizuno, Chihiro Shimizu, and  
Tsutomu Watanabe

This paper is available online at [www.gsb.columbia.edu/cjeb/research](http://www.gsb.columbia.edu/cjeb/research)

---

C O L U M B I A   U N I V E R S I T Y   I N   T H E   C I T Y   O F   N E W   Y O R K

# On the Evolution of the House Price Distribution

Takaaki Ohnishi\*    Takayuki Mizuno†    Chihiro Shimizu‡    Tsutomu Watanabe§

First draft: April 22, 2010

This version: May 27, 2011

## Abstract

Is the cross-sectional distribution of house prices close to a (log)normal distribution, as is often assumed in empirical studies on house price indexes? How does the distribution evolve over time? To address these questions, we investigate the cross-sectional distribution of house prices in the Greater Tokyo Area for the period 1986 to 2009. We find that size-adjusted house prices follow a lognormal distribution except for the period of the housing bubble and its collapse in Tokyo, for which the price distribution has a substantially heavier right tail than that of a lognormal distribution. In addition, we find that, during the bubble era, the sharp price movements were concentrated in particular areas, and this spatial heterogeneity is the source of the fat upper tail. These findings suggest that the shape of the size-adjusted price distribution, especially the shape of the tail part, may contain information useful for the detection of housing bubbles. Specifically, the presence of a bubble can be safely ruled out if recent price observations are found to follow a lognormal distribution. On the other hand, if there are many outliers, especially near the upper tail, this may indicate the presence of a bubble, since such price observations are unlikely to occur if they follow a lognormal distribution. This method of identifying bubbles is quite different from conventional ones based on aggregate measures of housing prices, and therefore should be a useful tool to supplement existing methods.

*JEL Classification Number:* R10; C16

*Keywords:* house price indexes; lognormal distributions; power-law distributions; fat tails; hedonic regression; housing bubbles; market segmentation

---

\*Correspondence: Takaaki Ohnishi, Canon Institute for Global Studies and University of Tokyo. E-mail: [ohnishi.takaaki@canon-igs.org](mailto:ohnishi.takaaki@canon-igs.org). We would like to thank Donald Haurin, Christian Hilber, Tomoyuki Nakajima, Kiyohiko G. Nishimura, Misako Takayasu, Hiroshi Yoshikawa, and the participants of the UNECE/ILO meeting on consumer price indexes in Geneva and Skye Seminar 2010 for helpful comments and suggestions. This research is a part of the project on “Understanding Inflation Dynamics of the Japanese Economy” funded by a JSPS Grant-in-Aid for Creative Scientific Research (18GS0101). Ohnishi gratefully acknowledges a grant from the Ministry of Education, Culture, Sports, Science, and Technology of Japan (Grant-in-Aid for Young Scientists (B) No. 20760053).

†University of Tsukuba. E-mail: [mizuno@cs.tsukuba.ac.jp](mailto:mizuno@cs.tsukuba.ac.jp)

‡Reitaku University. E-mail: [cshimizu@reitaku-u.ac.jp](mailto:cshimizu@reitaku-u.ac.jp)

§Hitotsubashi University and University of Tokyo. E-mail: [watanabe1284@gmail.com](mailto:watanabe1284@gmail.com)

# 1 Introduction

Researchers on house prices typically start their analysis by producing a time series of the *mean* of prices across different housing units in a particular region by, for example, running a hedonic or repeat-sales regression. In this paper, we pursue an alternative research strategy: we look at the entire distribution of house prices across housing units in a particular region at a particular point of time and then investigate the evolution of such cross-sectional distribution over time. We seek to describe price dynamics in the housing market not merely by changes in the mean but by changes in some key parameters that fully characterize the entire cross-sectional price distribution.

Specific questions we will address in this paper are as follows. First, we would like to know whether the price distribution is close to a normal distribution, as is often assumed in empirical studies on house price indexes, or whether it has fatter tails than a Gaussian distribution. Second, we are interested in how the shape of the price distribution is affected by house attributes, including the size and location of a house. Third, we would like to know how the shape of the distribution changes over time, especially during the housing bubble Japan experienced in the late 1980s and its burst in the early 1990s.

Recent studies on the cross-sectional distribution of house prices include Gyourko et al. (2006), McMillen (2008), Van Nieuwerburgh and Weill (2010), and Maattanen and Tervio (2010). The main interest of Gyourko et al. (2006), Van Nieuwerburgh and Weill (2010), and Maattanen and Tervio (2010) is the relationship between the house price distribution and the income distribution. For example, Maattanen and Tervio (2010) ask whether the recent increases in income inequality in the United States have had any impact on the distribution of house prices. On the other hand, McMillen (2008) focuses on the change in the house price distribution over time and asks whether the change in the price distribution comes from a change in the distribution of house characteristics such as size, location, age, and so on, or from a change in the implicit prices associated with those characteristics. The focus of our paper is closely related to the issues discussed in these papers, but differs from them in some important respects. First, this paper is the first attempt to specify the shape of the house price distribution, paying particular attention to the tail part of the distribution. Second, this paper examines the effect of a housing bubble on the cross-sectional price distribution. While steep rises of the *mean* of house prices in various countries in recent decades has received a lot of attention in the literature, the change in the shape of the cross-sectional price distribution has received much less attention. In this paper we seek to fill this gap.

Our main findings are as follows. First, the cross-sectional distribution of house prices has a fat upper tail and the tail part is close to that of a power law distribution. This is confirmed by the goodness-of-fit test recently proposed by Malevergne et al. (2009). On the other hand,

the cross-sectional distribution of house sizes, as measured by the floor space, has an upper tail that is less fat than that of the price distribution and is close to an exponential distribution. These two findings suggest a particular functional form of hedonic regression to identify the size effect. We construct a size-adjusted price by subtracting the house size (multiplied by a positive coefficient) from the log price and find that the size-adjusted price follows a lognormal distribution for most of the sample period. An important exception is the period of the housing bubble and its collapse in 1987-1995, during which the price distribution in each year has a power law tail even after controlling for the size effect.

Second, we divide the entire sample area into small pixels and find that the size-adjusted price is close to a lognormal distribution *within* each of these pixels even during the bubble period, but its mean and variance are highly dispersed *across* different pixels. This finding implies that the sharp price hike during the bubble period was concentrated in particular areas, and this spatial heterogeneity is the source of the fat upper tail observed for the bubble period.<sup>1</sup> We interpret this as evidence for market segmentation during the bubble period.

The rest of the paper is organized as follows. Section 2 explains the dataset and the empirical strategy we employ. Sections 3 and 4 present our size- and location-adjustments to house prices. Section 5 concludes the paper.

## 2 Data and Empirical Strategy

### 2.1 Data

We use a unique dataset that we have compiled from individual listings in a widely circulated real estate advertisement magazine, which is published on a weekly basis by Recruit Co., Ltd., one of the largest vendors of residential lettings information in Japan. The dataset covers the Greater Tokyo Area for the period 1986 to 2009, including the bubble period in the late 1980s and its collapse in the first half of the 1990s. It contains 724,416 listings for condominiums and 1,602,918 listings for single family houses.<sup>2</sup> In this paper we will use data only for condominiums. According to Shimizu et al. (2004), this dataset covers more than 95 percent of the entire transactions in the central part of Tokyo (namely, the 23 special wards of Tokyo), although its coverage for suburbs is limited. This dataset is used by a series of papers, including Shimizu et al. (2010), which compares hedonic and repeat-sales measures in terms

---

<sup>1</sup>Cochrane (2002) argues that an important feature of the tech stock bubble in the late 1990s is that it was concentrated in stocks related to internet business. Cochrane (2002: 17) states that “if there was a ‘bubble,’ or some behavioral overenthusiasm for stocks, it was concentrated on Nasdaq stocks, and Nasdaq tech and internet stocks in particular.”

<sup>2</sup>The dataset contains full information about the evolution of the posted price for a housing unit from the week when it is first listed until the final week when it is removed because of successful transaction. In this paper, we use the price only at the final week since it can be safely regarded as sufficiently close to the contract price. The number of listings shown in the text does not include those prices listed before the final week.

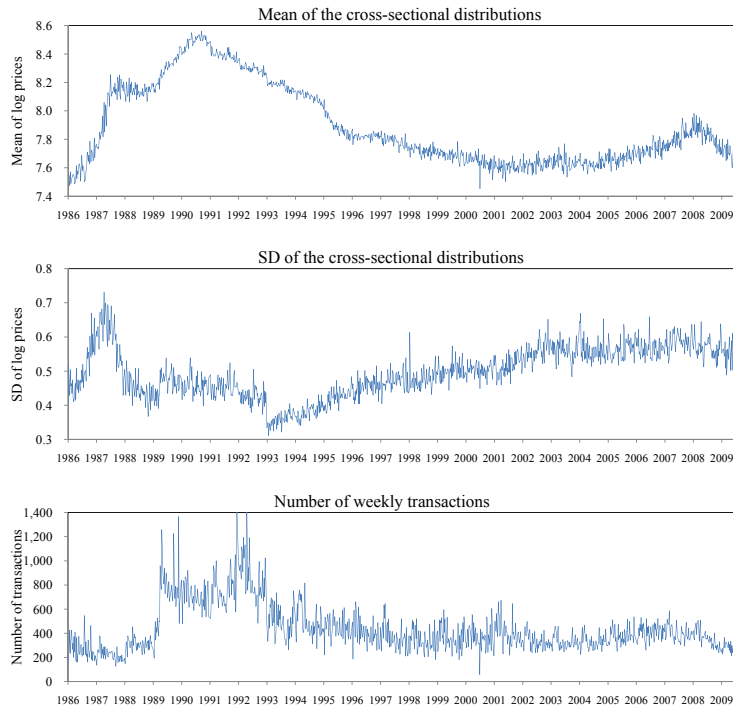


Figure 1: Weekly fluctuations in prices and transaction volume

of their performance.

Figure 1 shows changes in the mean of the cross-sectional house price distribution in the upper panel, the standard deviation in the middle panel, and the transaction volume in the lower panel. We see that the mean price exhibits a sharp increase between the beginning of 1987 and the beginning of 1988. Previous studies refer to this as the first phase of the housing bubble in Tokyo. After a short break in 1988, prices started to rise again in 1989 and continued to rise until the fall of 1990. This is the second phase of the housing bubble. Soon after the fall of 1990, prices started to turn down, followed by a slow but persistent decline for more than a decade until prices bottomed out in 2002, when the mean price reached the level before the bubble started in 1987. Prices finally began to rise again in 2003 and continued to rise until registering a sharp decline in 2008 due to the recent global financial crisis.

Turning to the standard deviation shown in the middle panel, this exhibits a sharp rise during the first phase of the bubble and stayed high during the second phase.<sup>3</sup> Finally, the

---

<sup>3</sup>We also see a secular increase in price dispersion since 1993. We are not quite sure why this is the case, but recent studies, including Van Nieuwerburgh and Weill (2010) and Gyourko et al. (2006) find some evidence that the recent rise in house price dispersion across regions in the United States is related to the change in income distribution across regions. For example, Van Nieuwerburgh and Weill (2010) find that the cross-sectional coefficient of variation (CV) of house prices across 330 metropolitan statistical areas in the United States increased from 0.15 in 1975 to 0.53 in 2007. Through a counterfactual simulation, they show that this

bottom panel, which shows the transaction volume, indicates that the number of transactions exhibits a sharp increase at the beginning of 1989, exactly when the mean price started to rise, although the transaction volume remained practically unchanged during the first phase of the bubble. Somewhat interestingly, the transaction volume remained at a high level even in 1991 and 1992, when the mean price had already started to decline.

## 2.2 Empirical strategy

A widely used approach to deal with product heterogeneity in terms of quality is hedonic analysis and there are numerous applications to housing services. The core idea of hedonic analysis is that the value of a product is the sum of the values of product characteristics. For example, Shimizu et al. (2010) start their analysis by assuming that the value of a house is the sum of the values of attributes such as its floor space, its age, the commuting time to the nearest station, and so on, and run hedonic regressions using these attributes as independent variables.

This idea has important implications regarding the shape of the cross-sectional distribution of house prices. To show this, let us start by assuming that the price of house  $i$  at a particular point in time, which is denoted by  $P_i$ ,<sup>4</sup> is the sum of  $K$  components:

$$P_i = F(X_{i1}, X_{i2}, \dots, X_{ik}, \dots, X_{iK}) \quad (1)$$

where  $P_i$  and  $X_{ik}$  are both random variables and  $X_{i1}, \dots, X_{iK}$  are assumed to be independent from each other. Furthermore, we assume a multiplicative functional form such that

$$P_i = \prod_{k=1}^K X_{ik}. \quad (2)$$

Taking logarithm of both sides of this equation leads to:

$$\ln P_i = \sum_{k=1}^K x_{ik} \quad (3)$$

where  $x_{ik} \equiv \ln X_{ik}$ . This equation appears frequently in hedonic analyses of house prices. It simply states that the price of a house is equal to the sum of  $K$  random variables.

Given this setting, the central limit theorem tells us that the sum of these random variables converges to a normal distribution if the number of attributes,  $K$ , goes to infinity. Let us denote the variance of  $x_{ik}$  by  $s_k^2$ , and define the average variance  $\bar{s}_K^2$  as

$$\bar{s}_K^2 \equiv \frac{1}{K} (s_1^2 + s_2^2 + \dots + s_K^2).$$

---

increase in dispersion of house prices is accounted for mostly by the increase in income inequality.

<sup>4</sup>Note that the subscript for time is dropped here to simplify the exposition.

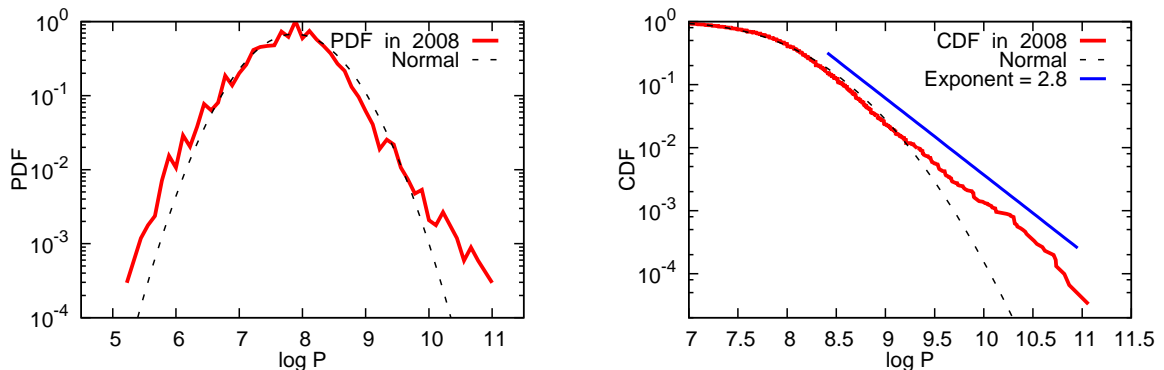


Figure 2: House price distribution in 2008

Then, according to the Lindberg-Feller central limit theorem, the sum of random variables  $\sum_{k=1}^K x_{ik}$  converges to a normal distribution as  $K$  goes to infinity if the average variance  $\bar{s}_K^2$  converges to a finite constant (namely,  $\lim_{K \rightarrow \infty} \bar{s}_K^2 = \bar{s}^2$ ) and the following condition is satisfied:

$$\lim_{K \rightarrow \infty} \frac{\max_{k \leq K} \{s_k\}}{K \bar{s}_K} = 0. \quad (4)$$

In other words, the theorem states that the sum of random variables, regardless of their form, will tend to be normally distributed. A notable feature of this result is that it does *not* require that the variables in the sum come from the same underlying distribution. Instead, the theorem requires only that no single term dominates the average variance, as stated in (4). Put differently, condition (4) states that none of the random variables is dominantly large relative to their sum.<sup>5</sup> A famous textbook example of the central limit theorem is the distribution of persons' height. The height distribution of, say, mature men of a certain age can be considered normal, because height can be seen as the sum of many small and independent effects. Similarly, the log price of houses will be normally distributed if house prices are determined as the sum of many small and independent effects.

The above argument suggests that the lognormal distribution can be seen as a benchmark for the cross-sectional distribution of house prices. However, some previous studies on house price distributions find that the actual distributions have fatter tails than a lognormal distribution. For example, McMillen (2008), using data on single family houses in Chicago for 1995, shows that the kernel density estimates for the log price are asymmetric, with a much fatter lower tail. Against this background, we examine the extent to which the house price distribution deviates from a lognormal distribution using our observations for 2008. The results are presented in Figure 2, where the left panel shows the probability density function (PDF), with the horizontal axis representing the yen price in logarithm and the vertical axis representing

<sup>5</sup>For more on this theorem, see Feller (1968). Greene (2003) provides a compact description of various versions of the central limit theorem including this one.

the corresponding density, also in logarithm. The empirical distribution is shown by the red line while the lognormal distribution with the same mean and standard deviation is shown by the black dotted line. The figure indicates that the tails of the empirical distribution are fatter than those of the lognormal distribution. In particular, the upper tail of the empirical distribution is much fatter than that of the lognormal distribution. To examine the differences in the upper tail more closely, we accumulate the densities from the right (upper) tail to produce the cumulative distribution function (CDF), which is shown in the right panel. In this panel, the value on the vertical axis corresponding to the value of 9.2 on the horizontal axis, for example, is 0.01, meaning that the fraction of houses whose prices are equal to or higher than that price level is 1 percent. We now see more clearly that the upper tail of the empirical distribution is fatter than that of the lognormal distribution. For example, the fraction of housing units whose price deviates from the mean by more than  $3\sigma$  is about 1.47 percent, while the corresponding number for the lognormal distribution is only 0.26 percent.

What causes the empirical distribution to deviate from the benchmark (i.e., the lognormal distribution)? This is the main topic we address in this paper. Our hypothesis is that some of the factors that determine house prices are dominantly volatile, so that condition (4) is violated. Denoting these dominant factors by vector  $Z_i$ , the house price distribution,  $\Pr(P_i = p)$ , can be decomposed as follows:

$$\Pr(P_i = p) = \sum_z \Pr(P_i = p \mid Z_i = z) \Pr(Z_i = z). \quad (5)$$

Note that the house price distribution conditional on  $Z_i$ , namely  $\Pr(P_i = p \mid Z_i = z)$ , should be a lognormal distribution, since the dominant factors are now fully controlled for. This means that the right-hand side of equation (5) is a weighted sum of lognormals, with the weights being given by  $\Pr(Z_i = z)$ . We know that the sum of lognormals with different means and variances is no longer a lognormal (see, for example, Feller (1968)), and the hypothesis we examine is that this is why the house price distribution deviates from the benchmark. Given this hypothesis, we proceed as follows in the remainder of the paper: we first specify the dominant factors and then eliminate them, thereby constructing prices that are adjusted for these factors; finally, we examine whether these adjusted prices follow a lognormal distribution.

Diewert et al. (2010) argue that there are three important price determining characteristics: the land area of the property; the livable floor space area of the structure; and the location of the property. Similarly, previous studies on house prices in Japan, including Shimizu et al. (2010), find that the floor space of a housing unit (especially in the case of condominiums) and its location play dominantly important roles in determining its price. This empirical evidence suggests that the size and the location of a property are key candidates for the  $Z$  variables. We will identify and eliminate the size effect in the next section, and the location effect in Section 4.



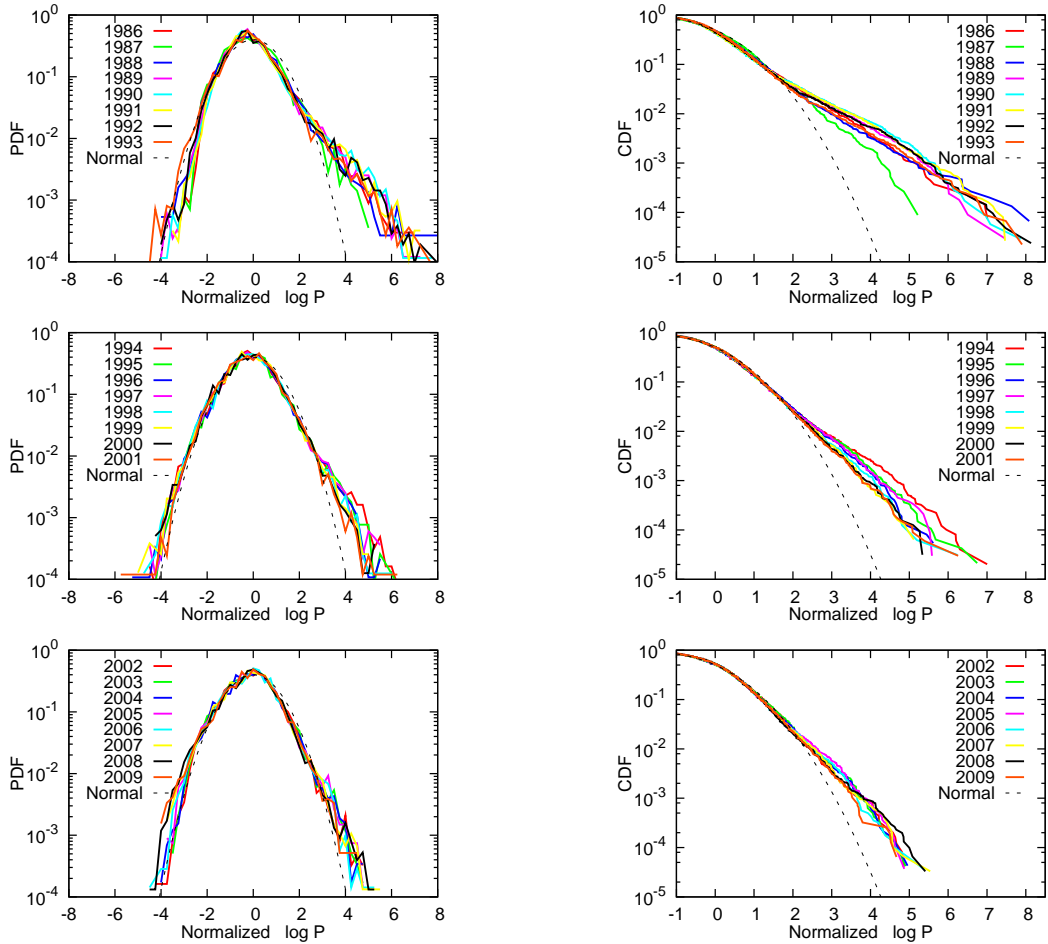


Figure 3: House price distributions by year

### 3 Size-adjustment to House Prices

#### 3.1 Distribution of unadjusted house prices

Figure 3 presents the PDFs and the CDFs of the cross-sectional price distribution for each year from 1986 to 2009. To make the price distributions in different years comparable, we normalize the log prices in year  $t$  by subtracting the mean in year  $t$  (i.e., the mean of log prices in year  $t$ ) and dividing by the standard deviation in year  $t$  (i.e., the standard deviation of log prices in year  $t$ ). The lognormal lines in the figure represent the CDF of a standard lognormal distribution. Note that the CDFs are constructed in the same way as in Figure 2, that is, the value on the vertical axis corresponding to a price level is the sum of the densities *above* that price level.

The first thing we see from the figure is that, as in Figure 2, the PDFs and the CDFs show fatter upper tails than a lognormal distribution. More importantly, we see that the deviation

from a lognormal distribution tends to be larger in the late 1980s and the first half of the 1990s. Specifically, the PDFs in these years are substantially skewed to the right, indicating that during the bubble period house prices did not rise by the same percentage for every housing unit; instead, price increases were concentrated in particular housing units, so that relative prices across houses changed significantly.

The CDFs in this figure provide more detailed information regarding the shape of the price distributions. We see that the CDF for each year forms an almost straight line in this log-log graph, implying that the house price distribution is well approximated by a power law distribution (or a Pareto distribution) at least in the tail part, the PDF and CDF of which are given by

$$\Pr(P_{it} = p) = \frac{\zeta_t m_t^{\zeta_t}}{p^{\zeta_t+1}}; \quad \Pr(P_{it} \geq p) = \left(\frac{m_t}{p}\right)^{\zeta_t}; \quad p > m_t > 0 \quad (6)$$

where  $P_{it}$  denotes the price of house  $i$  in period  $t$ , and  $\zeta_t$  and  $m_t$  are time-variant positive parameters.<sup>6</sup> The shape of a power law distribution is mainly determined by the parameter  $\zeta_t$ , which is referred to as the exponent of the power law distribution. Smaller values for  $\zeta_t$  imply fatter tails. Note that the CDF given in (6) implies that

$$\ln \Pr(P_{it} \geq p) = -\zeta_t \ln p + \zeta_t \ln m_t.$$

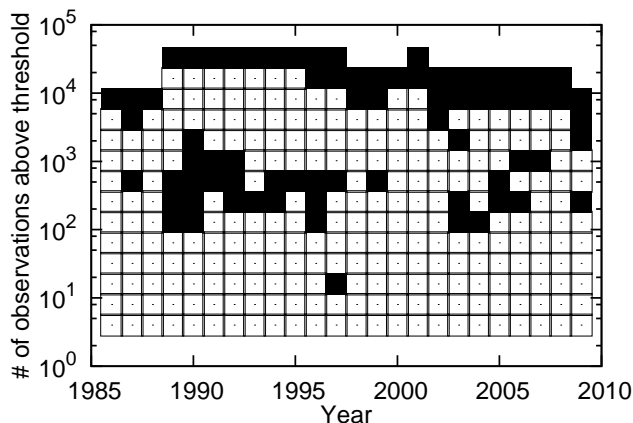
In other words, the log of the cumulative probability should be linearly related to the log price, and the slope of the linear line between the two variables should be equal to  $-\zeta_t$ . The CDFs in Figure 3 suggest the presence of such a linear relationship between the log price and the log of the cumulative probability. We see from the CDF in Figure 2 that  $\zeta_{2008}$  is about 2.8. Similarly, we find from the corresponding figures for the other years (which are not shown due to space limitations)  $\zeta$  also all take values of around three.<sup>7</sup>

As a goodness-of-fit test, we employ the test proposed by Malevergne et al. (2009). Specifically, we test the null hypothesis that, beyond some threshold  $u$ , the upper tail of the house price distribution is characterized by a power law distribution

---

<sup>6</sup>See Gabaix (2008) for an extensive survey of empirical and theoretical studies on power laws in various economic contexts such as income and wealth, the size of cities and firms, and stock market returns.

<sup>7</sup>Note that we cannot obtain estimates for  $\zeta_t$  from Figure 3. The CDFs in Figure 3 are for normalized prices, which are defined by  $[P_{it} \exp(-\mu_t)]^{1/\sigma_t}$ , where  $\mu_t$  and  $\sigma_t$  are the mean and the standard deviation in year  $t$ . Therefore, the slope of each CDF in Figure 3 is given by  $\sigma_t \zeta_t$  (rather than  $\zeta_t$ ) if the original price follows the power law distribution given by (6). Taleb (2007) provides many examples of power law distributions. For example, the net worth of Americans follows a power law distribution with an exponent of 1.1; the frequency of the use of words follows such a distribution with an exponent of 1.2; the population of U.S. cities has an exponent of 1.3; the number of hits on websites has an exponent of 1.4; the magnitude of earthquakes has an exponent of 2.8; and market moves have an exponent of 3 (or lower). The exponents for the house price distributions estimated here are greater than most of these figures, implying that the tail part of the house price distributions are less fat than in the other examples of power law distributions.



- : The null (power law) is rejected at the 1 percent significance level.
- : The null is not rejected at the 1 percent significance level.

Figure 4: Power law distribution versus lognormal distribution

$$\Pr(P = p; \alpha) = \alpha \cdot \frac{u^\alpha}{p^{\alpha+1}} \cdot 1_{p \geq u}$$

against the alternative that the upper tail follows a lognormal beyond the same threshold, i.e.,

$$\Pr(P = p; \alpha, \beta) = \left[ \sqrt{\frac{\pi}{\beta}} \exp\left(\frac{\alpha^2}{4\beta}\right) \left(1 - \Phi\left(\frac{\alpha}{\sqrt{2\beta}}\right)\right) \right]^{-1} \frac{1}{p} \exp\left(-\alpha \ln \frac{p}{u} - \beta \ln^2 \frac{p}{u}\right) \cdot 1_{p \geq u}$$

where  $\Phi(\cdot)$  represents the CDF of a standard normal distribution. Note that this is equivalent to testing the null that the upper tail of the *log* price follows an exponential distribution against the alternative that it follows a normal distribution. For this transformed test, Del Castillo and Puig (1999) have shown that the clipped empirical coefficient of variation  $\hat{c} \equiv \min(1, c)$  provides the uniformly most powerful unbiased test, where  $c$  is the empirical coefficient of variation. The result of our goodness-of-fit test is presented in Figure 4, where the horizontal axis represents the year and the vertical axis represents the number of observations above the threshold  $u$ . For example,  $10^3$  on the vertical axis means that the threshold  $u$  is set such that the number of observations above  $u$  is  $10^3$ . A black square indicates that the null is rejected at the 1 percent significance level for a particular year-threshold combination, while a white square indicates that the null is not rejected at the same significance level. The figure shows that a power law distribution provides a good approximation for the 500 most expensive houses, while a lognormal distribution provides a better approximation for the set of less expensive houses.

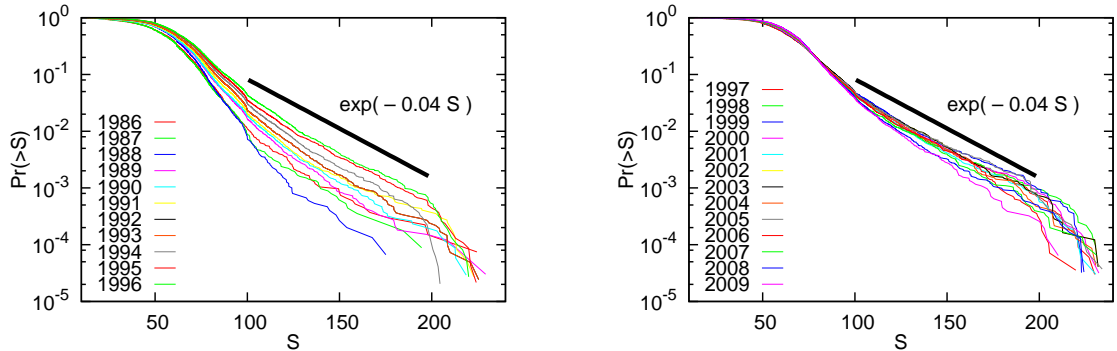


Figure 5: Cumulative house size distributions

### 3.2 Distribution of house sizes

Previous studies on wealth (or income) distributions across households have typically found that those distributions are characterized by fat upper tails, and that they follow a power law distribution (see Pareto (1896)). Given that houses form an important part of households' wealth, it may be not that surprising to detect a similar pattern in house price distributions. However, the result that house prices follow a power law distribution is not consistent with the argument based on the central limit theorem. Why do house prices follow a power law distribution rather than a lognormal distribution? As a first step to address this question, we decompose the house price distribution as follows:

$$\Pr(P_{it} = p) = \sum_s \Pr(P_{it} = p \mid S_i = s) \Pr(S_i = s), \quad (7)$$

where  $S_i$  represents the size of housing unit  $i$ , which is measured by the floor space of that unit. The term  $\Pr(S_i = s)$  represents the distribution of house sizes, while the term  $\sum \Pr(P_{it} = p \mid S_i = s)$  represents the distribution of house prices conditional on house size. An important thing to note is that even if each of these conditional distributions is lognormal, the weighted sum of lognormals with different mean and variance is not a lognormal distribution. This is a potential source of the power law tails that we observed in our house price data.

We start by examining the term  $\Pr(S_i = s)$  in equation (7). Figure 5 presents the CDFs of house sizes for each year, with the floor space, measured in square meters, on the horizontal axis and the log of the CDF on the vertical axis. We see that the CDF for each year is close to a straight line in this semi-log graph, implying that the size distribution can be approximated by an exponential distribution whose PDF and CDF are given by

$$\Pr(S_i = s) = \lambda_t \exp(-\lambda_t s); \quad \Pr(S_i \geq s) = \exp(-\lambda_t s); \quad \lambda_t > 0. \quad (8)$$

Note that the CDF shown above implies that

$$\ln \Pr(S_i \geq s) = -\lambda_t s,$$

so that the log of the CDF depends linearly on house size. This is what we see in Figure 5. The slope of the CDF line, namely the value of  $\lambda$ , is almost identical for the different years and is somewhere around 0.04.

The fact that house sizes follow an exponential distribution implies that the tails of the size distribution are less fat than those of the price distribution. For example, for 2008, the fraction of housing units whose size deviates from the mean by more than  $3\sigma$  is only 0.94 percent, while the corresponding number for the price distribution is 1.47 percent.<sup>8</sup>

### 3.3 Size-adjusted prices

We now turn to the relationship between the price of a house and its size, which is represented by the conditional probability  $\Pr(P_{it} = p \mid S_i = s)$  in equation (7). We propose a hedonic model which is consistent with the fact that house prices and sizes follow, respectively, a power law distribution with an exponent of  $\zeta_t$  and an exponential distribution with an exponent of  $\lambda_t$ . That is, the log prices are determined as

$$\ln P_{it} \sim \left( \frac{\lambda_t}{\zeta_t} \right) S_i + \epsilon_{it}, \quad (9)$$

where  $\epsilon_{it}$  is a normally distributed random variable, which, as we saw in Section 2.2, can be interpreted as the sum of many small and independent factors. To show equation (9), we first note that the PDF of the exponential distribution given in (8) implies that  $(\lambda_t/\zeta_t)S_i$  follows an exponential distribution with an exponent of  $\zeta_t$  if  $S_i$  itself is an exponential distribution with an exponent of  $\lambda_t$ . Next, we can show that the sum of the random variable that follows an exponential distribution and the random variable that follows a normal distribution is well approximated only by the exponential distribution when the sum takes large values (because of the much fatter tails of an exponential distribution).<sup>9</sup> Combining the two, the right-hand side of (9) is well approximated by an exponential distribution with an exponent of  $\zeta_t$  when the sum of the two terms on the right-hand side takes large values. On the other hand, the fact that  $P_{it}$  follows a power law distribution with an exponent of  $\zeta_t$  implies that  $\ln P_{it}$  follows an exponential distribution with an exponent of  $\zeta_t$ . In this way we can confirm that each side of equation (9) follows an identical distribution with an identical exponent.<sup>10</sup>

---

<sup>8</sup>To see why the tails of the house size distribution are less fat than the tails of the price distribution, consider a simple example in which household A has 100 times as much wealth as household B, so that A spends 100 times as much money on a house as B. What does A's house look like? Does it have a bathroom that is 100 times larger than the one in B's house? Alternatively, does it have 100 bathrooms? Needless to say, neither is true, because even a person of A's wealth would have little use for such a gigantic bathroom (or so many bathrooms). Instead, it is more likely that the size of A's house (and therefore the size and/or number of its bathroom) is only, say, 10 times greater and, consequently, the unit area price of A's house, 10 times higher than B's.

<sup>9</sup>See the appendix for a formal proof of this.

<sup>10</sup>The price-size relationship described by equation (9) provides an answer to the question regarding the choice of functional form for hedonic price equations, which has been extensively discussed by previous studies such

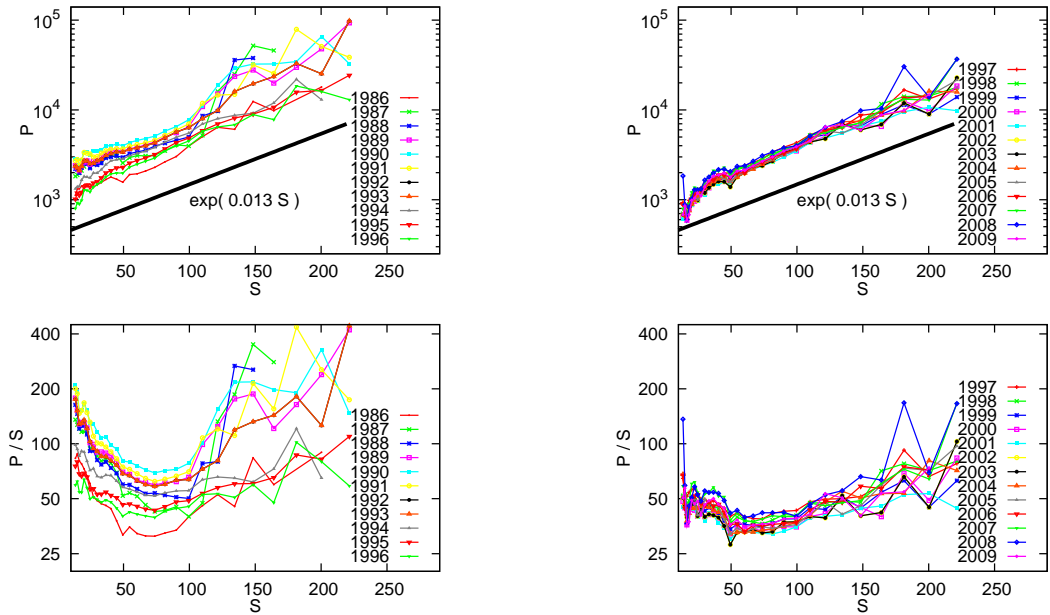


Figure 6: Relationship between house size and price

To empirically test the hedonic model given by (9), we first examine for a linear relationship between the log price of houses and their size. The upper panels of Figure 6 show the floor space on the horizontal axis and the median of the log price corresponding to that size on the vertical axis. These panels indicate that there exists a stable linear relationship between the two variables. Furthermore, equation (9) implies that the per unit area price,  $P/S = [\exp(\lambda/\zeta)S + \text{positive constant}]/S$ , decreases with  $S$  when  $S$  is small and increases with  $S$  when  $S$  is sufficiently large, so that there should exist a U-shaped relationship between the per unit area price and the house size. The lower panels of Figure 6, in which the vertical axis now measures  $P/S$ , confirms this prediction.

Second, we run an OLS regression of the form

$$\ln P_{it} = a_t S_i + b_t + \eta_{it} \quad (10)$$

to see whether the disturbance term  $\eta_{it}$  is indeed normally distributed as assumed in (9). The regression results are presented in Figures 7 and 8. Figure 7 shows the estimates of  $a$  and  $b$  for each year. The estimate of  $a$  is almost identical across years and is around 0.013, implying that an increase in the house size by a square meter leads to a 1.3 percent increase in the house price. More importantly, the estimate of  $a$  is very close to the value predicted by (9). That is, the value of  $\zeta$  is around 3 as we saw in Section 3.1, and the value of  $\lambda$  is about 0.04 as

---

as Cropper et al. (1988), Diewert (2003), and Triplett (2004). The novelty of our approach is that we derive this functional form not from economic theories but from the statistical fact that house prices and sizes follow a power law and an exponential distribution, respectively.

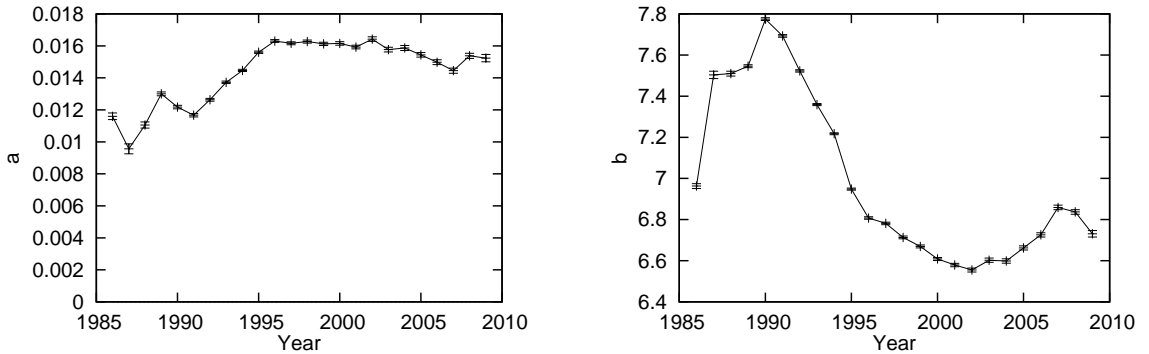


Figure 7: Price-size regressions

we saw in Section 3.2, so that the coefficient on  $S_i$ , namely  $\lambda/\zeta$ , should be something around 0.013 ( $= 0.04/3$ ). This is quite close to the point estimate of  $a$  for each year.<sup>11</sup> Turning to the estimate of  $b$ , this exhibits substantial fluctuations: it increases by more than 20 percent per year from 1986 to 1990 and then declines by 10 percent per year from 1990 to 2002.

Figure 8 shows the CDFs of size adjusted prices defined by

$$\tilde{P}_{it} \equiv \left[ P_{it} \exp(-\hat{a}_t S_i - \hat{b}_t) \right]^{1/\hat{\sigma}_t}, \quad (11)$$

where  $\hat{a}_t$  and  $\hat{b}_t$  are the estimates of  $a_t$  and  $b_t$ , and  $\hat{\sigma}_t$  is the estimate for the standard deviation of  $\eta_{it}$ . Note that the hedonic model given by (9) implies that  $\tilde{P}_{it}$  should be a lognormal distribution. The CDFs of the size adjusted prices are shown in the three panels on the right-hand side of Figure 8, while the price distributions *without* size adjustments (the same figures as in Figure 3) are shown on the left-hand side. Comparing these two sets of CDFs, we see that the CDFs of the size-adjusted prices are much closer to the CDF of a lognormal distribution. More specifically, the CDFs for 2002 to 2009, which are shown in the lower right panel, are almost identical to the CDF of a lognormal distribution. The same applies to the CDFs for 1996 to 2001, which are shown in the middle right panel. However, the CDFs for 1986 to 1995, which are presented in the upper right and the middle right panels, are still far from the CDF of a lognormal distribution, although they are slightly closer to it than the CDFs of the non-adjusted prices.

## 4 Location Adjustment to House Prices

The analysis in the previous section suggested that size-adjusted prices followed a lognormal distribution at least for quiet periods without large price fluctuations. This is consistent with

<sup>11</sup>Note that the per unit area price,  $\exp(aS + b)/S$  takes its minimum value when  $S$  is equal to  $1/a$ . Given the estimate of  $a$ , this implies that the per unit area price takes its minimum value when  $S = 1/0.013 \approx 75$ , which is consistent with what we see in the lower two panels of Figure 6.

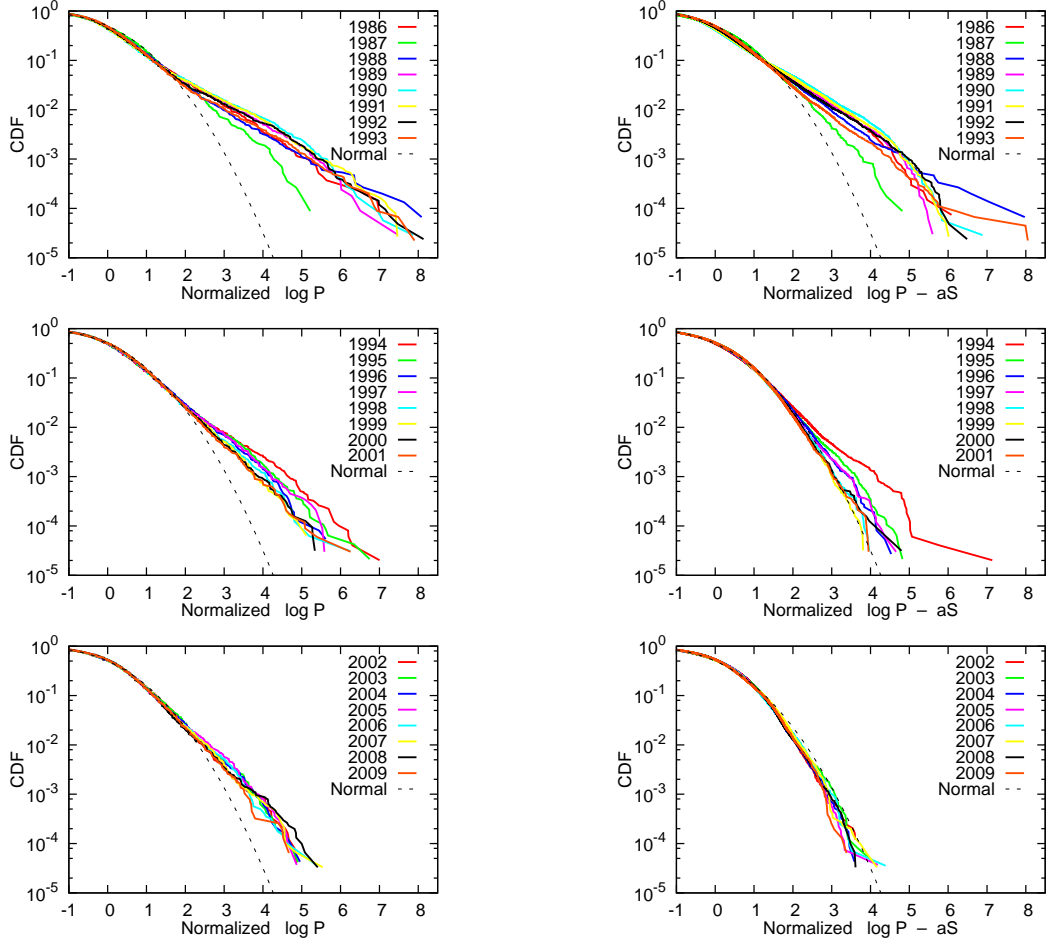


Figure 8: Cumulative distributions of size-adjusted house prices

the idea that, as stated in (7), the power law tails of the original prices stem from the mixture of lognormal distributions with different mean and variance. At the same time, the analysis in the previous section showed that the fat tails of the price distribution remain largely unchanged for the bubble period (i.e., the late 1980s and the first half of the 1990s) even after controlling for the size effect. This suggests that there still remains some mixture of lognormal distributions.

In this section, we test the hypothesis that the power law tails of the size adjusted prices during the bubble period arise due to the mixture of different lognormal distributions corresponding to different regions. To do so, we start by decomposing the size-adjusted price into the sum of conditional distributions:

$$\Pr\left(\tilde{P}_{i,r,t} = p\right) = \sum_{\theta} \Pr\left(\tilde{P}_{i,r,t} = p \mid \theta_{rt} = \theta\right) \Pr\left(\theta_{rt} = \theta\right) \quad (12)$$

where  $\tilde{P}_{i,r,t}$  denotes the size-adjusted price for a house located in region  $r$ , which is defined by



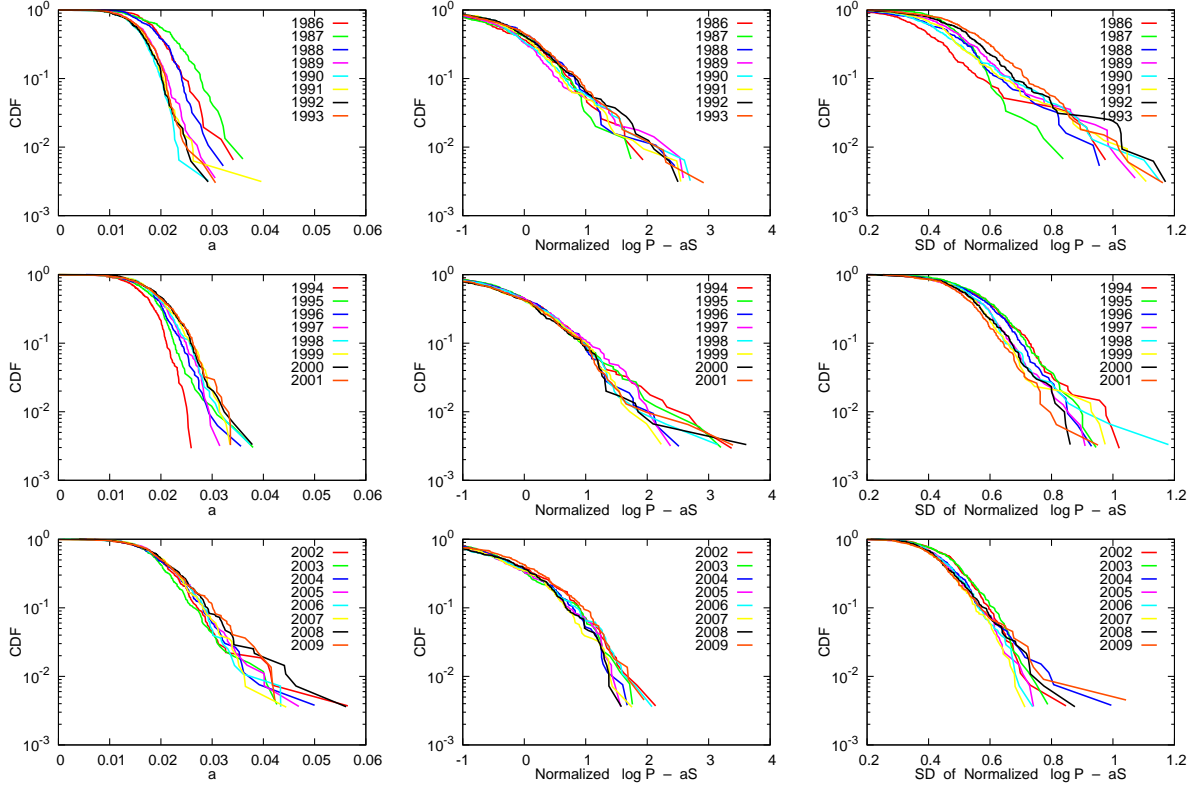


Figure 9: Dispersion of  $a_r$ ,  $b_r$  and  $\sigma_r$  across pixels

$\tilde{P}_{i,r,t} \equiv P_{i,r,t} \exp(-a_{rt}S_{i_r} - b_{rt})$ . The vector of parameters  $\theta_{rt}$  is defined by

$$\theta_{rt} \equiv (a_{rt}, b_{rt}, \sigma_{rt}), \quad (13)$$

where the parameters  $a_{rt}$ ,  $b_{rt}$ , and  $\sigma_{rt}$  are the coefficient on the house size variable, the constant term, and the standard deviation of the disturbance term in equation (10), but it is assumed in this section that they could differ depending on the location of a house. The location effect is fully controlled for in the conditional distributions  $\Pr(\tilde{P}_{i,r,t} = p \mid \theta_{rt} = \theta)$ , so that they should be lognormals. According to equation (12), the distribution of  $\tilde{P}_{i,r,t}$  is a mixture of these lognormals, each of which is for a different region.

We first examine the distribution of  $\theta_{rt}$  across different regions. Specifically, we divide the Greater Tokyo Area into pixels of 0.033 degrees latitude and 0.033 degrees longitude or roughly 3.3 by 3.3 km.<sup>12</sup> Then, using size-adjusted prices within a pixel, we run a regression of the form:

$$\ln P_{i,r,t} = a_{rt}S_{i_r} + b_{rt} + \eta_{i,r,t} \quad (14)$$

for each combination of  $r$  and  $t$  and obtain  $\hat{\theta}_{rt} \equiv (\hat{a}_{rt}, \hat{b}_{rt}, \hat{\sigma}_{rt})$ . The regression results are

<sup>12</sup>Note that one degree is approximately 100 km.

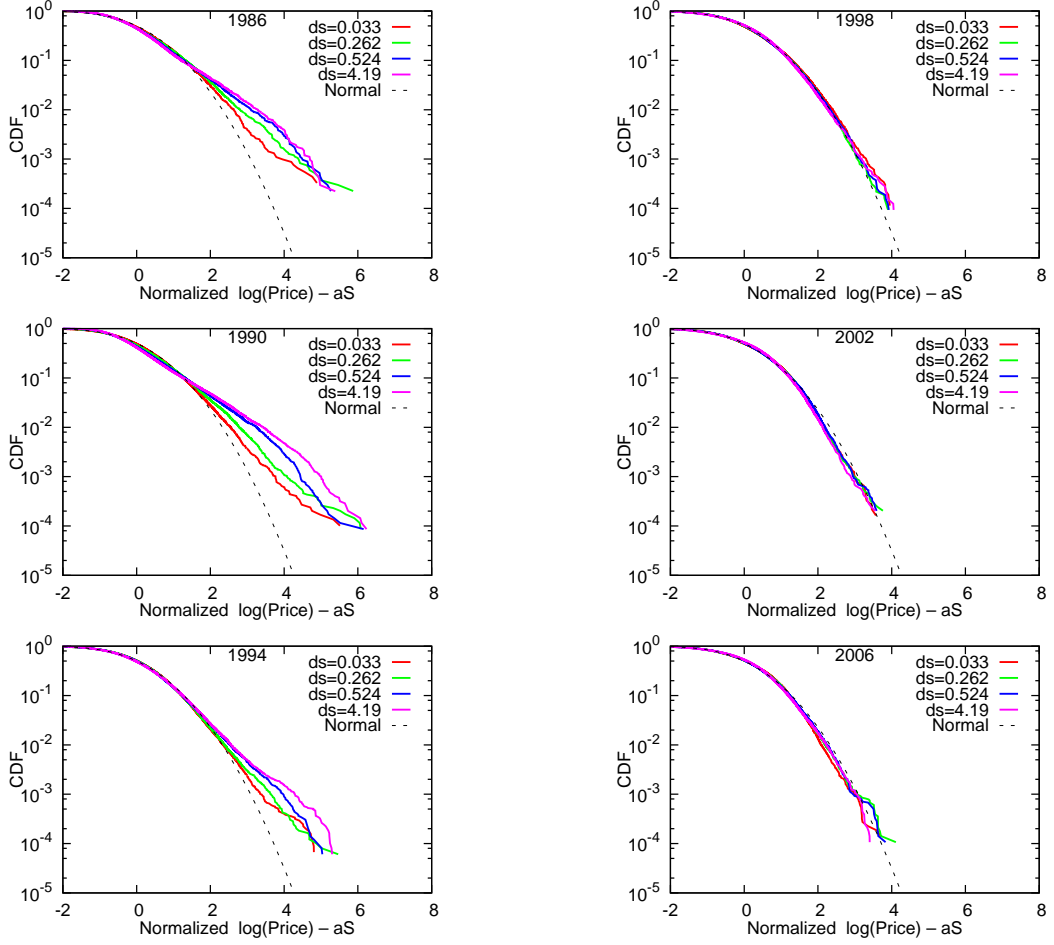


Figure 10: Cumulative distributions of size-adjusted house prices for different pixel sizes

presented in Figure 9.<sup>13</sup> The three left-most panels show the CDFs of  $\hat{a}_{rt}$ , while the panels in the middle and the right-most panels respectively show the CDFs of  $\hat{b}_{rt}$  and  $\hat{\sigma}_{rt}$ . The CDFs of  $\hat{a}_{rt}$  indicate that  $a$  is less dispersed across pixels during the period of the bubble and its collapse (1987-1995) than in the other years. On the other hand, the CDFs of  $\hat{b}_{rt}$  and  $\hat{\sigma}_{rt}$  show that these parameters are more highly dispersed during the same period, implying that the sharp price hike during the bubble period was concentrated in particular pixels. Put differently, the housing market was segmented during this period.

Next, we investigate whether the conditional distributions are close to a lognormal distribution. Using the estimates of  $\theta_{rt}$  obtained from the regression, we calculate the size-adjusted prices for each pixel:

$$\tilde{P}_{i,t} \equiv \left[ P_{i,t} \exp(-\hat{a}_{rt} S_{i,t} - \hat{b}_{rt}) \right]^{1/\hat{\sigma}_{rt}}. \quad (15)$$

<sup>13</sup>In conducting these regressions, we use only those pixels with more than twenty transactions in a year. The number of pixels used in the regressions is about 300 for each year.

The estimated CDFs of  $\tilde{P}_{i,t}$  are presented in Figure 10 for the years 1986, 1990, 1994, 1998, 2002, and 2006. Note that each of the six panels contains four different lines, each of which corresponds to a different pixel size, namely 4.190 by 4.190 degrees, 0.524 by 0.524 degrees, 0.263 by 0.263 degrees, and 0.033 by 0.033 degrees. The results for 1998, 2002, and 2006 indicate that the CDFs are very close to a lognormal distribution, irrespective of pixel size. This is not very surprising given that, as we saw in the previous section, the CDFs in these years were already close to a lognormal distribution before controlling for the location effect. For the period of the bubble and its collapse, we see more interesting results: for 1986, 1990, and 1994, the estimated CDF tends to become closer to a lognormal distribution as the pixel size becomes smaller.<sup>14</sup>

In sum, the analysis in this section shows that the distribution of size-adjusted prices *within* a pixel is fairly close to a lognormal distribution even during the period of the bubble and its collapse, but its mean and standard deviation are highly dispersed *across* different pixels. As a result, the sum of these lognormals turns out to be far from a lognormal distribution during this period. In other words, heterogeneity across pixels in terms of mean and standard deviation is the source of the fat upper tail of the size-adjusted price distribution during the period of the bubble and its collapse.

## 5 Summary and Some Policy Implications

In this paper, we found that the cross-sectional distribution of house prices in the Greater Tokyo Area has a fat upper tail, and the tail part is close to that of a power law distribution. On the other hand, the cross-sectional distribution of house sizes measured in terms of floor space has less fat tails than the price distribution and is close to an exponential distribution. We proposed a hedonic model consistent with these findings and confirmed that size-adjusted prices follow a lognormal distribution except for the period of the asset bubble and its collapse in Tokyo for which the price distribution remains asymmetric and skewed to the right even after controlling for the size effect. As for the period of the bubble and its collapse, we found some evidence that the sharp price movements were concentrated in particular areas, and this spatial heterogeneity is the source of the fat upper tail.

The analysis in this paper shows that the cross-sectional distribution of size-adjusted prices is very close to a lognormal distribution during regular times but deviated substantially from a lognormal during the bubble period. This suggests that the shape of the size-adjusted price distribution, especially the shape of the tail part, may contain information useful for the de-

---

<sup>14</sup>It should be noted that the estimated CDF does not fully converge to a lognormal even in the case of the smallest pixels. It may be the case that the CDF becomes much closer still to a lognormal distribution if we were able to reduce the pixel size even further. Unfortunately, we cannot do so because of the limited number of observations.

tection of housing bubbles. That is, the presence of a bubble can be safely ruled out if recent price observations are found to follow a lognormal distribution. On the other hand, if there are many outliers, especially near the upper tail, this may indicate the presence of a bubble, since such price observations are very unlikely to occur if they follow a lognormal distribution. This method of identifying bubbles is quite different from conventional ones based on aggregate measures of housing prices, which are estimated either by hedonic or repeat-sales regressions, and therefore should be a useful tool to supplement existing methods.

## References

- [1] Cochrane, J. H. (2002), "Stocks as Money: Convenience Yield and the Tech-Stock Bubble," *NBER Working Paper* 8987.
- [2] Cropper, M. L., L. B. Deck, and K. E. McConnell (1988), "On the Choice of Functional Form for Hedonic Price Functions," *Review of Economics and Statistics* 70(4): 668-675.
- [3] Del Castillo, J., and P. Puig (1999), "The Best Test of Exponentiality Against Singly Truncated Normal Alternatives," *Journal of the American Statistical Association* 94, 529-532.
- [4] Diewert, W. E. (2003), "Hedonic Regressions: A Consumer Theory Approach," in R. C. Feenstra and M. D. Shapiro (eds.), *Scanner Data and Price Indexes*, National Bureau of Economic Research Studies in Income and Wealth, Vol. 64, Chicago, IL: University of Chicago Press, 317-48.
- [5] Diewert, W. E., J. de Haan, and R. Hendriks (2010), "The Decomposition of a House Price Index into Land and Structures Components: A Hedonic Regression Approach," *Discussion Paper* 10-01, University of British Columbia.
- [6] Feller, W. (1968), *An Introduction to Probability Theory and Its Applications*. Vol. 1, Third Edition, New York: Wiley.
- [7] Gabaix, X. (2008), "Power Laws in Economics and Finance," *NBER Working Paper* 14299.
- [8] Green, W. (2003), *Econometric Analysis*, Fifth Edition, Upper Saddle River, NJ: Prentice Hall.
- [9] Gyourko, J., C. Mayer, and T. Sinai (2006), "Superstar Cities," *NBER Working Paper* 12355.
- [10] Maattanen, N., and M. Tervio (2010), "Income Distribution and Housing Prices: An Assignment Model Approach," *CEPR Discussion Paper* 7945.

- [11] Malevergne, Y., V. Pisarenko, and D. Sornette (2009), “Gibrat’s Law for Cities: Uniformly Most Powerful Unbiased Test of the Pareto against the Lognormal,” *American Economic Review*, forthcoming.
- [12] McMillen, D. P. (2008), “Changes in the Distribution of House Prices over Time: Structural Characteristics, Neighborhood, or Coefficients?” *Journal of Urban Economics* 64(3): 573-589.
- [13] Shimizu, C., K. G. Nishimura, and Y. Asami (2004), “Search and Vacancy Costs in the Tokyo Housing Market: An Attempt to Measure Social Costs of Imperfect Information,” *Review of Urban and Regional Development Studies* 16(3): 210-230.
- [14] Shimizu, C., K. G. Nishimura, and T. Watanabe (2010), “House Prices in Tokyo: A Comparison of Repeat-Sales and Hedonic Measures,” *Research Center for Price Dynamics Discussion Paper* 62, Hitotsubashi University.
- [15] Taleb, N. N. (2007), *The Black Swan: The Impact of the Highly Improbable*, New York: Random House, Inc.
- [16] Triplett, J. (2004), “Handbook on Hedonic Indexes and Quality Adjustments in Price Indexes: Special Application to Information Technology Products,” *OECD Science, Technology and Industry Working Papers* 2004/9.
- [17] Van Nieuwerburgh, S., and P. O. Weill (2010), “Why Has House Price Dispersion Gone Up?” *Review of Economic Studies*, forthcoming.