1 *Article*

# 2 A Machine Learning Approach to the Residential
# 3 Relocation Distance of Households Living in the
# 4 Seoul Metropolitan Region

5 **Changhyo Yi [1,*] and Kijung Kim [2]**

6 [1] Department of Urban Engineering, Hanbat National University, Daejeon 34158, Republic of Korea;
7 yich@hanbat.ac.kr
8 [2] Department of Urban Planning and Design, University of Seoul, Seoul 02504, Republic of Korea;
9 kimkj87@uos.ac.kr
10 * Correspondence: yich@hanbat.ac.kr; Tel.: +82-42-821-1194

11

12 **Abstract:** This study aimed to ascertain the applicability of a machine learning approach to the
13 description of residential mobility patterns of households in the Seoul metropolitan region (SMR).
14 The spatial range and temporal scope of the empirical study were set to 2015 to review the most
15 recent residential mobility patterns in the SMR. The analysis data used in this study involve the
16 microdata of Internal Migration Statistics provided by the Microdata Integrated Service of Statistics
17 Korea. We analysed the residential relocation distance of households in the SMR by using machine
18 learning techniques such as ordinary least squares regression and decision tree regression. The
19 results of this study showed that a decision tree model can be more advantageous than ordinary
20 least squares regression in terms of the explanatory power and estimation of moving distance. A
21 large number of residential movements are mainly related to the accessibility to employment
22 markets and some household characteristics. The shortest movements occur when households with
23 two or more members move into densely populated districts. In contrast, job-based residential
24 movements have relatively longer distance. Furthermore, we derived knowledge on residential
25 relocation distance, which can provide significant information on the urban management of
26 metropolitan residential districts and the construction of reasonable housing policies.

27 **Keywords:** residential relocation distance; residential movement; machine learning; decision tree
28 regression; Seoul metropolitan region

29

## 30 1. Introduction

31     A large number of households experience multiple residential movements during their lifetime,
32 although some people continue their lives in only one location. Residential movements have been
33 researched in residential choices and preferences as a searching process of appropriate location and
34 dwelling with respect to individual characteristics. However, residential choices and preferences
35 should be clearly distinguished. Residential choices indicate the actual behaviour related to
36 residential movement, and residential preference is related to the relative attractiveness of housing
37 and residential environment that affect movers [1]. Residential mobility can be represented by the
38 spatial moving pattern based on actual behaviours of the movers. Previous studies of spatial patterns
39 of residential relocations focused on the conventional research topics: frequency, direction, and
40 distance of residential mobility. The life-cycle model [2], sector model [3], and Ravenstein's Laws [4]
41 are representative research achievements, as well as theories related to these subjects. However, the
42 relevant empirical studies for household units have paid relatively scant attention to the topics of
43 direction and distance of residential mobility. This could be because of the excessive complexity of
44 influencing factors, lack of computing power to handle large volumes of data on household
45 movements, and absence of appropriate analytical model [5].

46    After economic achievement and quantitative growth [6], the Korean housing market has
47 experienced structural changes in terms of both supply and demand. The housing shortage problem
48 of the Korean society is considered to have been resolved with the housing supply ratio exceeding
49 100% in the early 21st century, and a fundamental change in the nature of household [2], which is a
50 basic unit of residential mobility and location change, is in progress. Representative phenomena
51 involve the reduction in the household size and aging, indicating the emergence of a new demand
52 class and the change of characteristics in the core demand groups. These situations are summarised
53 as the transitioning from a supply-based housing market to a demand-driven market [7]. Regarding
54 the demand, with the slowdown in population growth, the flow management of residential mobility,
55 considering the relocations within the metropolitan region, is gaining more importance than the
56 response to new demand caused by the increased population in the metropolitan region. Previous
57 studies [8,9] confirmed that the frequency of residential relocations of the Korean households is
58 relatively high among the Organization for Economic Cooperation and Development (OECD)
59 countries. In addition, recent studies [10,11] have shown that residential relocation distance could be
60 differentiated by household size and age of householder in the Seoul metropolitan region (SMR),
61 which is the most representative and largest metropolitan region in Korea. These phenomena could
62 be changed by the reduction of household size and aging trends.

63    An empirical understanding of spatial patterns and characteristics related to residential
64 relocation is important for the establishment of an in-depth housing policy. In addition, considering
65 the growing socio-economic complexity, residential mobility research using spatial Big Data is more
66 advantageous than research using only aggregated data. The study based on actual moving data of
67 households is more meaningful, as it can identify the practical residential moving patterns,
68 considering the conditions of a household rather than the ideal pursuit of a specific household. In a
69 continually changing housing market, such as the Korean housing market, the outcomes of such a
70 study could be applied to build a simulation model for forecasting future residential relocations.
71 Accordingly, academic reviews and empirical studies must attempt to apply new analytical methods
72 such as machine learning, which is used to derive meaningful knowledge from Big Data in the
73 housing and residential research fields. If such an attempt is successful, in the long run, it can be used
74 to construct a sustainable-housing-market-management system from the socio-economic aspect.

75    In this context, this study aimed to ascertain the applicability of a machine learning approach to
76 the description of the residential mobility patterns of households in the SMR. In particular, this study
77 focuses on the residential relocation distances of households, which is one of the main topics
78 representing residential spatial patterns and has not been focused upon in previous empirical studies
79 on household units, because residential relocation distance is a key factor in determining the spatial
80 extent of the housing (sub)market. In this paper, we first review literature on patterns and influencing
81 factors of residential relocation and examined the relocation characteristics of the SMR in Korea. Next,
82 we conducted empirical studies analysing the determinants in residential relocation distance by
83 using a machine learning approach. Finally, we conclude by summarising the outcomes of this study
84 and ascertaining the applicability of the machine learning method in estimating or forecasting studies
85 of housing and residential research.

## 2. Literature Review

87    Residential mobility is defined as a process of adjusting location to better meet the needs and
88 demands of a household [7,12,13]. Residential mobility can be divided into residential relocation and
89 urban migration. Residential relocation implies moving a residence within an urban living region,
90 and urban migration refers to moving out of that region. While urban migration mostly results from
91 changes in urbanisation and industrial structure [14], residential relocation is influenced by internal
92 and external factors of a household, such as income, composition, housing preference, and residential
93 environments. Residential movement occurs based on not only dissatisfaction with current location,
94 but also attractiveness to the new location [15–17]. Previous studies, which examined the influencing
95 factors of residential mobility, assumed a household-based decision-making mechanism. These
96 representative studies considered various household characteristics, such as composition of

household members [18], age and income [19], education level [20], and marriage duration [21]. However, these studies mainly focused on analysing the residential mobility.

Recently, not only the amount of flow but also the residential mobility patterns have gained interest in terms of suggesting implications to spatial planning and housing policy [11,22–25]. The moving patterns of households can be explained using the frequency, direction, and distance of residential mobility. In terms of the frequency of residential movements, the main reasons of residential mobility are the characteristics of the household and the changes in the life cycle of the household. In particular, the life cycle is a series of processes that human beings experience in their lives, resulting in a change in needs and demands for the living space according to each stage [2,26,27]. According to the life-cycle model, the changes in the characteristics of the frequency of movements depend on family events, such as marriage (formation), birth of children (expansion), moving out (contraction), and divorce or death of a spouse (dissolution) [2]. As the characteristics of the household change according to the life-cycle stage, many researchers have studied the probability of residential mobility affected by a stage. The previous studies show various empirical results in consideration of birth, childcare, marital age, and income with respect to individual households [18–21,26,28,29].

In terms of residential mobility direction, Hoyt's sector theory, which states "High grade residential growth tends to proceed from the given point of origin, along established lines of travel or toward another existing nucleus of buildings or trading centers" [3], is the initial theory in this research field. This theory suggests that the direction of residential mobility is due to the difference in rents generated in urban space. In the empirical study related to this theory, Burnley et al. [30] found that most of the residential mobilities in Australia are biased toward the outward direction from an urban center. Furthermore, Yang [31] reported that 26% of households moved to the outskirts of the city from the urban center, while only 9% of the households moved in the opposite direction. Regarding the distance of residential movement, the widely known Ravenstein's Laws suggests that most migration occurs over short distances [4]. The main research topics covered by related studies were concentrated to the quantity of flow of residential movements between origin and destination based on the gravity model. That is, the results of previous studies show a lack of in-depth research on the spatial patterns of residential moving distance. The short distance of residential movements is related to the existence of local housing markets (or housing submarkets) [32]. This study is a basic model that explains the residential relocation distance and links residential mobility to the local housing market. However, these previous studies did not consider demographic and socio-economic changes of modern society. In addition, while the studies on the frequency of residential movement considered the various characteristics of households, some studies on residential relocation distance and moving direction only considered the household characteristics.

Several studies have determined that the residential relocation distance differs according to household size, home ownership, job change, and parental status [10,11,29,33]. However, these studies compared and analysed the residential moving data aggregated by household characteristic. In addition, the model for estimating the moving distance of each household has not been developed yet. This is mostly because of the difficulty in obtaining the data of moved households and the lack of an analytical method for large volume data [5]. Nowadays, a large amount of residential relocation data of individual households is being provided by the Korean government agency and various analysing methods are being developed for Big Data. Especially in the Korean housing market, which is experiencing a rapid demographic change, the understanding of the spatial patterns of residential movements is gaining increasing importance because the housing demand and the behaviour of housing movement gradually change based on the household type. Therefore, this study focused on the application of a new approach that uses machine learning, which is advantageous for Big Data analysis, in order to empirically identify the impact of the household attributes and the location characteristics on the residential relocation distance in Korea.

146   **3. Characteristics of Residential Relocation Distance in SMR**

147       The main spatial range of this study was the SMR, which is a representative metropolitan region
148   located in northwestern South Korea and includes the cities of Seoul and Incheon and the Gyeonggi
149   province, and the temporal scope was the year 2015. The spatial unit of the present empirical analysis
150   involved the administrative district (Eup, Myeon, and Dong), which is a minimum-sized
151   administrative area-level unit in the SMR. The total area of the SMR is 11,828 km², with a population
152   of 23.906 million people living in 9.519 million households. In addition, the SMR contains two
153   metropolitan cities (Seoul and Incheon), one province (Gyeonggi-do), 28 cities (Si), 5 counties (Gun),
154   and 53 boroughs (Gu), which comprise 1,133 small administrative areas (Eup, Myeon, and Dong) (see
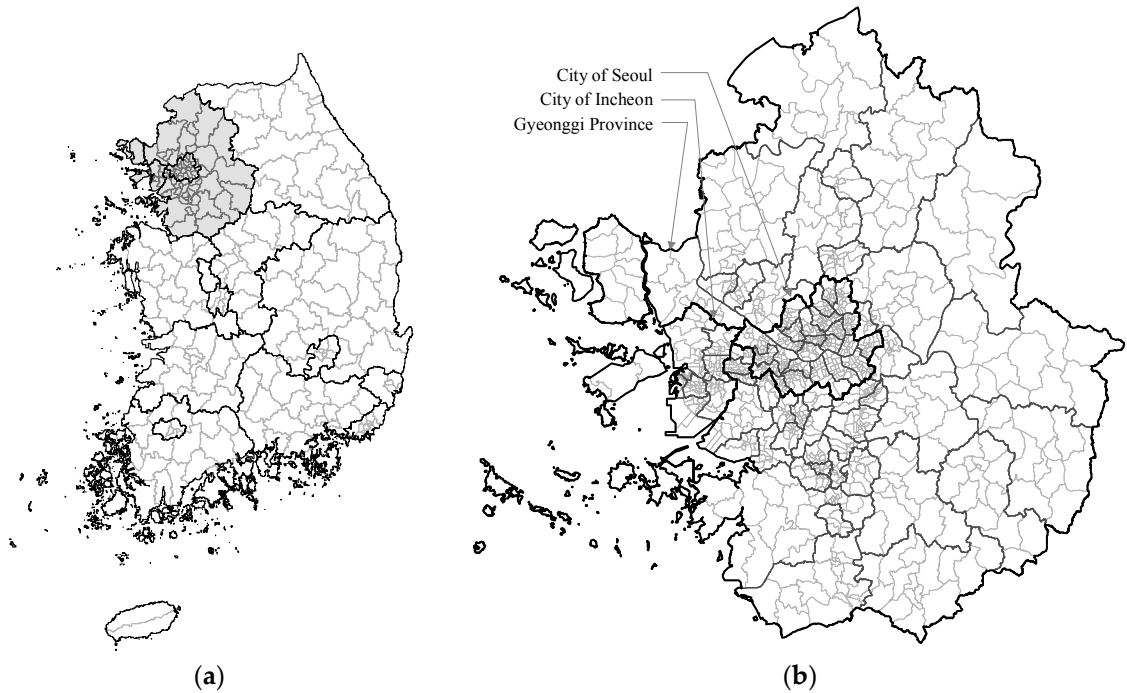155   Figure 1 and Table 1).

156



(**a**)                                    (**b**)

157                     **Figure 1.** (**a**) Location of SMR in Korea; (**b**) Components of SMR.

158                          **Table 1.** Seoul metropolitan region characteristics

| Item | | Total | Seoul | Incheon | Gyeonggi |
|---|---|---|---|---|---|
| Population (million people) | | 23.906 | 9.395 | 2.767 | 11.744 |
| Household (million households) | | 9.519 | 3.915 | 1.066 | 4.538 |
| Area (km²) | | 11,828 | 605 | 1,048 | 10,175 |
| City, county and borough level | Si | 28 | - | - | 28 |
| | Gun | 5 | - | 2 | 3 |
| | Gu | 53 | 25 | 8 | 20 |
| Minimum-sized administrative area level | Eup | 34 | - | 1 | 33 |
| | Myeon | 127 | - | 19 | 108 |
| | Dong | 972 | 424 | 129 | 419 |

159                 Source: Statistics of Urban Planning in 2015, 2015 Census in Korea.

160       The microdata of Internal Migration Statistics of Korea were used to analyse the spatial
161   characteristics of residential relocation. Internal Migration Statistics includes information of Korean
162   migrants from/to the smallest administrative areas of Eup, Myeon, and Dong obtained through using
163   the migrant's moving-in notifications. First, in the data analysis collected in 2015, the total number of
164   residential movements of households in Korea exceeds 6 million (6,098,915), of which approximately
165   3.1 million occurred in the SMR. The share of residential relocations within the SMR was 88.4%, which

166  occupied the majority of residential mobility in the metropolitan region. The share of residential
167  mobility in the metropolitan region was differentiated from the movement toward the inside and
168  outside by the municipality. The rates of residential relocations within the area were relatively low
169  in the metropolitan cities, such as Seoul and Incheon, and approximately 30% of residential
170  movements were confirmed to move beyond the boundaries of each municipality. The number of
171  residential movements per household was 0.326 in 2015, and the difference by area was not
172  significant. Second, the average residential relocation distance was 9.123 km in the SMR. As expected,
173  the average distance of residential movements from Seoul was the shortest (7.753 km), and that from
174  Gyeonggi province was the longest (10.391). However, the moving-out beyond the boundary of
175  Incheon city with the longest distance (29.112 km) was an unexpected outcome. This result is
176  presumed to be caused by the difference between the characteristics of the moving-out households
177  (refer to Table 2).

178

**Table 2.** Frequency and distance of residential movements

| Item | | SMR | Seoul | Incheon | Gyeonggi |
|---|---|---|---|---|---|
| Frequency of residential movements | Total | 3,107,134 (100.0%) | 1,287,379 (100.0%) | 352,488 (100.0%) | 1,467,267 (100.0%) |
| | Inside | 2,747,380 (88.4%) | 882,299 (68.5%) | 247,760 (70.3%) | 1,081,897 (73.7%) |
| | Outside | 359,754 (11.6%) | 405,080 (31.5%) | 104,728 (29.7%) | 385,370 (26.3%) |
| | Movement per household | 0.326 | 0.329 | 0.331 | 0.323 |
| Residential relocation distance [1] (km) | Total | 9.123 | 7.753 | 8.894 | 10.391 |
| | Inside | - | 3.940 | 4.304 | 7.965 |
| | Outside | - | 23.909 | 29.112 | 25.412 |

179    [1] The average Euclidian distance calculated using 10% randomly sampled data from the raw data.

180



(a)                                                                          (b)
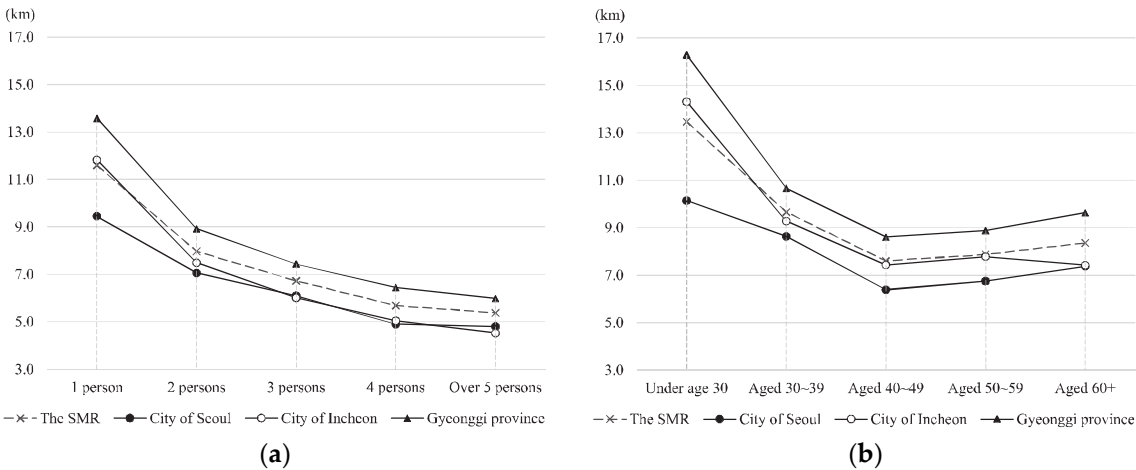
181    **Figure 2.** Relocation Distance based on (**a**) Family Size and (**b**) Age Group of Householder

182    Figure 2 shows the difference among residential moving distances by household types according
183  to characteristics of the household. In terms of household size, households with more members
184  represented a shorter moving distance. The households with three or more people in the metropolitan
185  cities (Seoul and Incheon) moved a similar distance, while the relocation distance of one-person
186  households showed a significant difference among municipalities. In addition, the age of a
187  householder is considered as a critical factor affecting the residential relocation distance of
188  households, and this result is identical with previously confirmed outcomes. The longest relocation
189  distance of households occurs for householders under age 30, and decreases to age range 40-49. Then,

190  the distance of residential movements gradually increases with age. This phenomenon agrees with
191  the results of previous literature [10,11] in Korea.
192      The estimated results of residential relocation distance in the SMR has the following implications.
193  First, the moving distance with respect to a household could vary according to the area in which the
194  household is located. Second, depending on the characteristics of the household members, there
195  could be differences in the moving distance. These outcomes imply that characteristics of households
196  and their location features should be considered in the construction of a model for estimating
197  residential relocation distance.

## 4. Methods and Materials

### 4.1. Decision Tree Using Machine Learning

200      The main analytical methodology of this study is machine learning, which is an efficient tool in
201  automatically detecting patterns of data and extracting information from large datasets [34]. Machine
202  learning differs from conventional statistics in that it is more concerned with making estimations or
203  predictions by using a model and formulating the generalisation process as a search through
204  hypotheses. In contrast, conventional statistics is more concerned with testing hypotheses [35].
205  Machine learning focuses on estimation or prediction by considering an optimal model, while the
206  latter concentrates on understanding the relationships between data. Recently, a few related studies
207  applying a machine learning-based method have been reported in various research fields, such as
208  environmental science, geomatics, and social science [36–39].
209      Decision trees in machine learning techniques are widely used for classification or regression
210  problems and generate the result in a tree form, which can be interpreted relatively easily compared
211  to the results of other techniques [40,41]. Thus, decision trees are known as a white-box model in the
212  software engineering field. Decision trees are classified into classification and regression trees, which
213  are constructed by repeatedly splitting the data. Each branch of a regression tree is partitioned
214  according to the homogeneity of two resulting groups; the homogeneity is maximised according to
215  the response variable. This method does not assume the relationship between the response and
216  predictors, unlike the conventional statistical model in which the relationship independent and
217  dependent variables is predefined and verified [42]. Therefore, the decision-tree regression method
218  has more advantages than the conventional statistical models with respect to fitting and estimation
219  using extremely complex data and structure. Therefore, in this study, the residential relocation
220  distance of each household in the SMR was analysed using decision-tree regression, which can be
221  regarded as the most appropriate model for analysis and estimation, considering rapidly changing
222  demographic transitions and household characteristics in the Korean housing market.

### 4.2. Selection of Explanatory Variables and Generation of Analysing Data

224      The estimated residential relocation distance is the dependent variable for conducting empirical
225  analysis by using a decision tree. The microdata obtained from the Internal Migration Statistics in
226  this study provides information of the smallest administrative district (Eup, Myeon, and Dong),
227  which is the same as a small-sized traffic analysis zone (TAZ), for the point of departure and
228  destination of each household's residential movement. Therefore, we estimated the moving distance
229  between the departure and destination based on the administrative center points by applying the
230  Euclidian distance calculation method. The cases for which the point of departure and destination
231  are the same, the following formula was applied to estimate the moving distance:

$$A = \pi r^2 \iff r = \sqrt{A/\pi}, \qquad\qquad (1)$$

232  where $A$ is the area of the administrative district and $r$ is the radius that assumes an irregularly
233  shaped administrative district as a circle.
234      The explanatory variables affecting the residential relocation distance of households moving
235  within the SMR were selected based on the results of previous studies and the hypothesis of the
236  present study. In this empirical analysis, not only the household attributes, but also the location

237 characteristics were selected considering the results from previous related research, for example, life-
238 cycle stages, residential mobility, and residential location choices. The variables contained in the
239 household attributes group were available from the microdata of Internal Migration Statistics.
240     In Table 3, the explanatory variables are classified into household attributes and location
241 characteristics. First, variables related to the attributes of household are *moving reason*, which include
242 job, house, and education; *age*; *gender*; *members*; *elderly people*; *children*; and *proportion of males* in the
243 household; these were collected from the microdata of Internal Migration Statistics in 2015. The three
244 nominal variables labeled as *moving reason* were coded as 1 if each moving reason was job, house, or
245 education, and 0 otherwise. These variables were selected to identify the influence of specific mobility
246 reasons of households on the moving distance. *Age* is defined as the age of the householder. *Gender*
247 is a nominal variable equal to 1 if the householder is male and 0 otherwise. *Member*, *elderly people*, and
248 *children* are variables related to the household structure; these are measured as the number of
249 corresponding members of each household. Finally, *proportion of male* is defined as the proportion of
250 male among total household members; it is measured at a ratio.

251     **Table 3.** Explanatory Variables Applied in the Empirical Analysis

| | Variable | Description | Unit | Source |
|---|---|---|---|---|
| Household attributes | Moving reason | Major reasons for residential relocation; Job/House/Education | - | The microdata of Internal Migration Statistics |
| | Age | Age of householder | Year | |
| | Gender | Male and female | - | |
| | Members | Number of household members | People | |
| | Elderly people | Number of elderly household members | People | |
| | Children | Number of school-aged children; Primary/Secondary | People | |
| | Proportion of males | Proportion of male household members | % | |
| Location characteristics [1] | Accessibility | Accessibility to employment market | - | Census on Establishments |
| | Density | Population density | People/ha | Population Census |
| | New building | Proportion of new building; 1 year/5 years | % | Housing Census |
| | Housing ownership | Ratio of owner-occupied housing | % | Population Census |
| | Rail availability | Ratio of rail catchment area | % | Korea Transport Database |
| | Bus availability | Number of metropolitan bus routes | EA | |

252    [1] The variables contained in the domain of location characteristics were calculated for both origin and
253                                       destination locations.

254     Second, the location variables include *accessibility*, *density*, *new building*, *housing ownership*, *rail*
255 *availability*, and *bus availability*, which were calculated with respect to both the departure and
256 destination positions of each household's residential movement. *Accessibility* was selected as an
257 explanatory variable measuring how the location advantage of employment opportunities affects the
258 moving distance of households. The accessibility to employment market was calculated using the
259 methodology representing location attraction, as mentioned by Hansen [43] and Wilson [44]:

$$Acc_i = \ln \sum_j Job_j \times \alpha\left(d_{ij}^{\beta}\right) \times exp\left(\gamma d_{ij}\right), \tag{2}$$

260 where $Acc_i$ is the accessibility of administrative district $i$; $Job_j$ is the number of jobs in potential
261 destination administrative district $j$; $d_{ij}$ represents the Euclidian distance between administrative
262 districts $i$ and $j$; and $\alpha$, $\beta$, $\gamma$ are the parameters. The parameters obtained from the analysis of
263 commuting patterns in the SMR in 2015 (the Metropolitan Transport Association) were applied in the
264 empirical analysis: 0.421 ($\alpha$), 0.276 ($\beta$), −0.082 ($\gamma$). *Density* is defined as the population density based
265 on administrative district, and *new building* is represented by the proportion of new buildings, that
266 is, the ratio of buildings that were constructed within the past 1 year (or 5 years). *Housing ownership*

267 is defined as the ratio of owner-occupied housing. These explanatory variables were selected to reflect
268 the influence of residential environments and housing conditions on the relocation distance of
269 households. In addition, two variables related to the availability of metropolitan transportation were
270 selected in this study. *Rail availability* is represented by the ratio of the catchment area within 500 m
271 from the metropolitan railway stations, and *bus availability* is defined as the number of metropolitan
272 bus routes operating in each administrative district.

273 *4.3. Descriptive Statistics*

274 This empirical analysis contains 209,252 residential movement data samples which is randomly
275 sampled 10% of the raw data including the householder information. The descriptive statistics for the
276 selected and estimated variables are listed in Table 4.
277 In the dataset, the average moving distance of households is 9.12 km, and the range of distance
278 is 0.24–267.31 km. Regarding the household attributes, 19% of the entire residential movements were
279 caused by a job. In addition, 60% and 2% of the residential relocations were due to housing
280 replacement and educational environment, respectively. Although the reasons for the residential
281 movements were numerous, housing replacement accounted for more than half. The average age of
282 householders was approximately 44.32, and the dataset consisted of 66% male and 34% female
283 population. The number of household members ranged from 1 to 9, with an average value of 2.1. On
284 average, the households included 0.14 elderly people, 0.12 primary school-aged children, and 0.14
285 secondary school-aged children. The proportion of males among household members was 53%.

286
**Table 4.** Descriptive Statistics

|  | Variable | Unit | Average | SD | Minimum | Maximum |
|---|---|---|---|---|---|---|
| - | Relocation distance | km | 9.12 | 13.66 | 0.24 | 267.31 |
| Household attributes | Moving reason: Job [1] | - | 0.19 | 0.39 | 0.00 | 1.00 |
|  | Moving reason: House [1] | - | 0.60 | 0.49 | 0.00 | 1.00 |
|  | Moving reason: Education [1] | - | 0.02 | 0.14 | 0.00 | 1.00 |
|  | Age | Years | 44.32 | 13.79 | 0.00 | 103.00 |
|  | Gender: Male [1] | - | 0.66 | 0.47 | 0.00 | 1.00 |
|  | Members | People | 2.10 | 1.30 | 1.00 | 9.00 |
|  | Elderly people | People | 0.14 | 0.41 | 0.00 | 4.00 |
|  | Children: Primary | People | 0.12 | 0.40 | 0.00 | 4.00 |
|  | Children: Secondary | People | 0.14 | 0.42 | 0.00 | 7.00 |
|  | Proportion of males | % | 53.52 | 38.27 | 0.00 | 100.00 |
| Location characteristics in origin | Accessibility | - | 14.26 | 0.53 | 6.40 | 14.79 |
|  | Density | People/ha | 174.25 | 129.14 | 0.00 | 550.00 |
|  | New building: 1 year | % | 2.93 | 2.46 | 0.27 | 17.36 |
|  | New building: 5 years | % | 13.69 | 6.12 | 1.83 | 34.89 |
|  | Housing ownership | % | 48.55 | 8.45 | 29.38 | 79.26 |
|  | Rail availability | % | 25.39 | 27.76 | 0.00 | 100.00 |
|  | Bus availability | EA | 7.54 | 10.11 | 0.00 | 71.00 |
| Location characteristics in destination | Accessibility | - | 14.23 | 0.54 | 6.40 | 14.79 |
|  | Density | People/ha | 167.34 | 129.29 | 0.00 | 550.00 |
|  | New building: 1 year | % | 3.11 | 2.74 | 0.27 | 17.36 |
|  | New building: 5 years | % | 14.04 | 6.26 | 1.83 | 34.89 |
|  | Housing ownership | % | 48.81 | 8.39 | 29.38 | 79.26 |
|  | Rail availability | % | 24.46 | 27.58 | 0.00 | 100.00 |
|  | Bus availability | EA | 7.66 | 10.28 | 0.00 | 71.00 |

287 [1] A reference of nominal variables.

288 As the residential relocation of household has a departure point and an arrival point, the location
289 characteristics were classified into not only the origin, but also the destination domains. As the

290  location characteristics of administrative districts are assigned to individual households, the
291  minimum and maximum values of characteristics at the origin and destination are the same. In
292  contrast, the differences in the averages and standard deviations are due to the number of households
293  included in each administrative district. The average values of accessibilities to origin and destination
294  were 14.26 and 14.23, respectively. In addition, the population density at the origin location (174.25
295  people/ha) was larger than that of the destination (167.34 people/ha). These results indicate that the
296  households moved out to less densely-populated districts. At the origin location, the proportion of
297  newly-constructed buildings within a year was 2.93% and that within five years was 13.69%.
298  Moreover, at the destination location, the proportion of newly constructed buildings within a year
299  was 3.11% and that within five years was 14.04%. These outcomes imply that the households moved
300  out to the districts with more new buildings in 2015. The ratio of rail catchment area in the origin
301  districts was 25.39% on average, which is larger than that in the destination districts (24.46%).
302  Moreover, the average number of bus routes was 7.66 at the destination and 7.54 at the origin location.

## 5. Results and Discussion

304      The analytical dataset was composed of 209,252 samples of residential households that moved
305  in 2015. In a machine learning approach, the analytical dataset is randomly split into training and
306  testing subsets. Generally, the former consists of 75% of the entire dataset and the latter consists of
307  the remaining 25%.

*5.1. Comparison of the Empirical Results between Ordinary Least Squares and Decision Tree Regressions*

309      In this study, the empirical analysis on residential relocation distance in the SMR included the
310  application of ordinary least squares regression and decision tree regression using a machine learning
311  approach. The results of the empirical analysis are summarised in Table 5.
312      First, the training and test R-squared values in the ordinary least squares regression model were
313  0.180 and 0.190, respectively, showing low explanatory power. In the household attributes domain,
314  among the residential moving reasons, *house* was an influencing factor that shortened the moving
315  distance of households by about 2 km compared to other reasons. This can be interpreted as a result
316  of the existence and influence of the housing sub-market in the SMR. On the other hand, *job* and
317  *education* were significant factors—these were significant factors affecting residential mobility in
318  previous studies [7,45]—in increasing the distance of residential movement of households over 5 km
319  compared to other causes. *Age* and *squared age* were significant variables, and the residential
320  relocation distance of households was the minimum at the householder age of approximately 59,
321  which is similar to the residential mobility of the life-cycle model. For the explanatory variables that
322  represent composition of a household, the number of household members and the number of children
323  had negative coefficients at the 99% level. These outcomes are similar to previous literature related
324  to mobility based on residential duration [46], which can be understood that households with more
325  members have more complex decision-making system for their residential relocation and there is a
326  tendency to maintain their community that was formed in the previous location. Whereas, *gender* was
327  a positive determinant at the 99% level. This result can be interpreted as the relatively low resistance
328  to residential moving distance in the households with a male householder or the long-distance
329  residential movements due to changes in the workplace of the male householder.
330      In the location characteristics domain, the most important explanatory variable was *accessibility*
331  to employment markets in the both the origin and destination residential locations. *Accessibility*
332  variables had negative coefficients, which implies the importance of proximity to employment
333  centers affecting residential location choice of household in previous studies [7,47–49]. *Density* and
334  proportions of *new buildings* within one year or five years also had significant coefficients, but their
335  signs showed opposite values in origin and destination locations of the residential movements. High
336  population density is considered as a negative determinant in residential environment [50], whereas
337  newly constructed houses are seen as a positive one. Since the former and the latter are a push factor
338  and a pull factor [51], respectively, the difference of the distance, as well as the migration flow of
339  intra-urban residential mobility can be generated. *Housing ownership* had negative coefficients in both

340  the origin and destination locations. These results can be interpreted as a relatively short movement
341  of residents living in the stabilised settlements based on the high proportion of housing ownership.
342  Moreover, the coefficient of *bus availability*, which is the number of inter-regional bus routes by
343  administrative district, showed a significant positive sign in only destination residential location. This
344  outcome means that even though it is located far away, the district with a large number of bus routes
345  with relatively high inter-regional mobility can be a residential moving destination.

346  **Table 5.** Results of the Empirical Analysis Using Machine Learning Models

| Variable (Feature) | | | Ordinary Least Squares Regression | | | | Decision Tree Regression | |
|---|---|---|---|---|---|---|---|---|
| | | | β | Std. β | Sig. | | Importance | Rank |
| | | (Constant) | 136.3587 | | 0.000 | ** | | |
| Household attributes | X(0) | Moving reason: Job | 5.6836 | 2.2401 | 0.000 | ** | 0.13180 | 3 |
| | X(1) | Moving reason: House | -1.9648 | -0.9630 | 0.000 | ** | - | - |
| | X(2) | Moving reason: Education | 5.4827 | 0.7688 | 0.000 | ** | 0.00289 | 8 |
| | X(3) | Age | -0.2362 | -3.2602 | 0.000 | ** | 0.00114 | 9 |
| | X(4) | Squared Age | 0.0020 | 2.7813 | 0.000 | ** | - | - |
| | X(5) | Gender: Male | 0.5935 | 0.2809 | 0.000 | ** | - | - |
| | X(6) | Members | -1.0025 | -1.3052 | 0.000 | ** | 0.01246 | 6 |
| | X(7) | Elderly people | 0.1631 | 0.0668 | 0.140 | | - | - |
| | X(8) | Children: Primary | -0.5795 | -0.2338 | 0.000 | ** | - | - |
| | X(9) | Children: Secondary | -0.8717 | -0.3697 | 0.000 | ** | - | - |
| | X(10) | Proportion of males | 0.0019 | 0.0739 | 0.154 | | - | - |
| Location characteristics in origin | X(11) | Accessibility | -2.0368 | -1.0766 | 0.000 | ** | 0.57976 | 1 |
| | X(12) | Density | 0.0008 | 0.1081 | 0.015 | * | 0.01450 | 5 |
| | X(13) | New building: 1 year | -0.1787 | -0.4411 | 0.000 | ** | - | - |
| | X(14) | New building: 5 years | -0.0461 | -0.2824 | 0.000 | ** | 0.00434 | 7 |
| | X(15) | Housing ownership | -0.0417 | -0.3523 | 0.000 | ** | 0.00001 | 12 |
| | X(16) | Rail availability | 0.0014 | 0.0392 | 0.352 | | - | - |
| | X(17) | Bus availability | 0.0055 | 0.0553 | 0.138 | | - | - |
| Location characteristics in destination | X(18) | Accessibility | -6.1138 | -3.2953 | 0.000 | ** | 0.23433 | 2 |
| | X(19) | Density | -0.0026 | -0.3358 | 0.000 | ** | 0.01749 | 4 |
| | X(20) | New building: 1 year | 0.1147 | 0.3156 | 0.000 | ** | - | - |
| | X(21) | New building: 5 years | 0.0434 | 0.2719 | 0.000 | ** | - | - |
| | X(22) | Housing ownership | -0.0271 | -0.2277 | 0.000 | ** | 0.00039 | 11 |
| | X(23) | Rail availability | 0.0019 | 0.0526 | 0.219 | | - | - |
| | X(24) | Bus availability | 0.0434 | 0.4457 | 0.000 | ** | 0.00090 | 10 |
| Explanatory Power | | | Training R²: 0.180 Test R²: 0.190 | | | | Training R²: 0.512 Test R²: 0.504 | |

347          * p-value < .05; ** p-value < .01.

348       Second, in decision trees, the complex tree constructed using the training dataset generally has
349  an overfitting problem. Therefore, by setting the parameters for maximum depth and the leaf node
350  minimum sample value, an early stopping method was applied to terminate the learning algorithm
351  before the tree became too complex [52]. The application of early stopping has advantages of not only
352  mitigating the overfitting problem, but also interpreting the derived tree structure. A trial and error
353  method was applied to set the appropriate parameter values: the maximum depth is 6, and the leaf
354  node minimum sample value is 10 (refer to Appendix A).
355       In the model applying decision tree regression, the explanatory powers of the final derived
356  model showed a remarkable improvement over the ordinary least squares regression model. The
357  training R-squared value was 0.512, and the test one was 0.504. Twelve features were contained in
358  the derived decision tree. The importance of features reflects the contribution each variable makes in

estimating the target variable, which is the residential relocation distance of each household in this study. The importance of a feature is estimated as the normalised total reduction of the criterion caused by the feature. In Table 5, two of the most important features were *accessibility* to employment markets in the locations before and after the residential movement. Among the residential moving reasons, *job* was ranked as the third most important feature. The importance of these three features accounted for approximately 95% of the total importance. In addition, in terms of importance, the following features were ranked: *density* of population in destination and origin locations, *members*, *new buildings* within five years in origin residential location, *moving reason: education*, *age* of householder, *bus availability* in destination, and *housing ownership* in destination and origin residential locations. These importance values are different from the standardised beta values that indicate the relative influence of explanatory variables on the results of the ordinary least squares regression model.

Decision trees, while not as powerful from a pure machine learning standpoint, are still one of the canonical examples of an understandable machine learning algorithm. That is, the structure of the derived decision tree can be represented, as shown in Figure 3. In the figure, the gray circle indicates a leaf node which is composed of 57 nodes, and intermediate nodes are represented using 56 white circles. Among them, the leftmost white circle is called the root node. In this study, the derived decision tree structure can be traced back to the splits from the training dataset starting with 156,939 samples at the root node. Moreover, in the tree structure, the solid lines mean that an observation goes to the lower branch if the condition shown at the intermediate node is satisfied, whereas the broken lines indicate that an observation goes to the upper branch if it is not satisfied. The equation presented on the right side of the intermediate node is a condition splitting the assigned samples of each node. Of the two numbers located to the right of the leaf nodes, the first and second numbers are the number of samples and the average relocation distance of the assigned samples to the nodes.

The black solid and broken lines represent branch paths that can reach the leaf nodes to which the top three most-allocated samples are assigned. In a decision tree model, describing the entire tree structure is not only extremely complex, but also inefficient. Therefore, the most important top three assigned leaf nodes and their assigned paths are described in this paper. First, the leaf node with the largest number of samples contains 43,880 households (27.96%) with an average residential moving distance of 6.866 km. The features affecting the path of branches to the leaf node were residential mobility caused by factors other than job or education ($X(0)$ and $X(2)$), higher potential accessibility to employment markets from origin to destination residential location ($X(11)$ and $X(18)$), and one-person household ($X(6)$). This result can be summarised as the pattern of general residential movements based on the employment market in the one-person household group. Second, the leaf node with the second largest number of samples includes 38,003 households (24.22%), with an average residential relocation distance of 3.312 km, which is the shortest distance leaf node in this derived tree. The features related to the leaf node were residential mobility caused by factors other than job ($X(0)$), higher potential accessibility to employment markets from origin to destination residential location ($X(11)$ and $X(18)$), densely populated origin and destination locations ($X(12)$ and $X(19)$), and households with more than two people ($X(6)$). This path can be understood as the shortest residential moving pattern of households with more than two members based on accessibility to employment markets, which are carried out among densely populated districts, for purposes other than a job. Finally, the path related to the third largest leaf node includes 9,732 households (6.20%) with an average moving distance of 8.323 km, which were affected by the features including residential mobility caused by job ($X(0)$) and lower potential accessibility to employment markets from origin to destination location ($X(11)$ and $X(18)$). This result can be interpreted as the relatively longer residential moving distance of households caused by job or employment except other residential conditions. In addition, Figure 3 describes that there are a lot of paths based on the decision tree of households related to the residential relocation distance in the SMR.
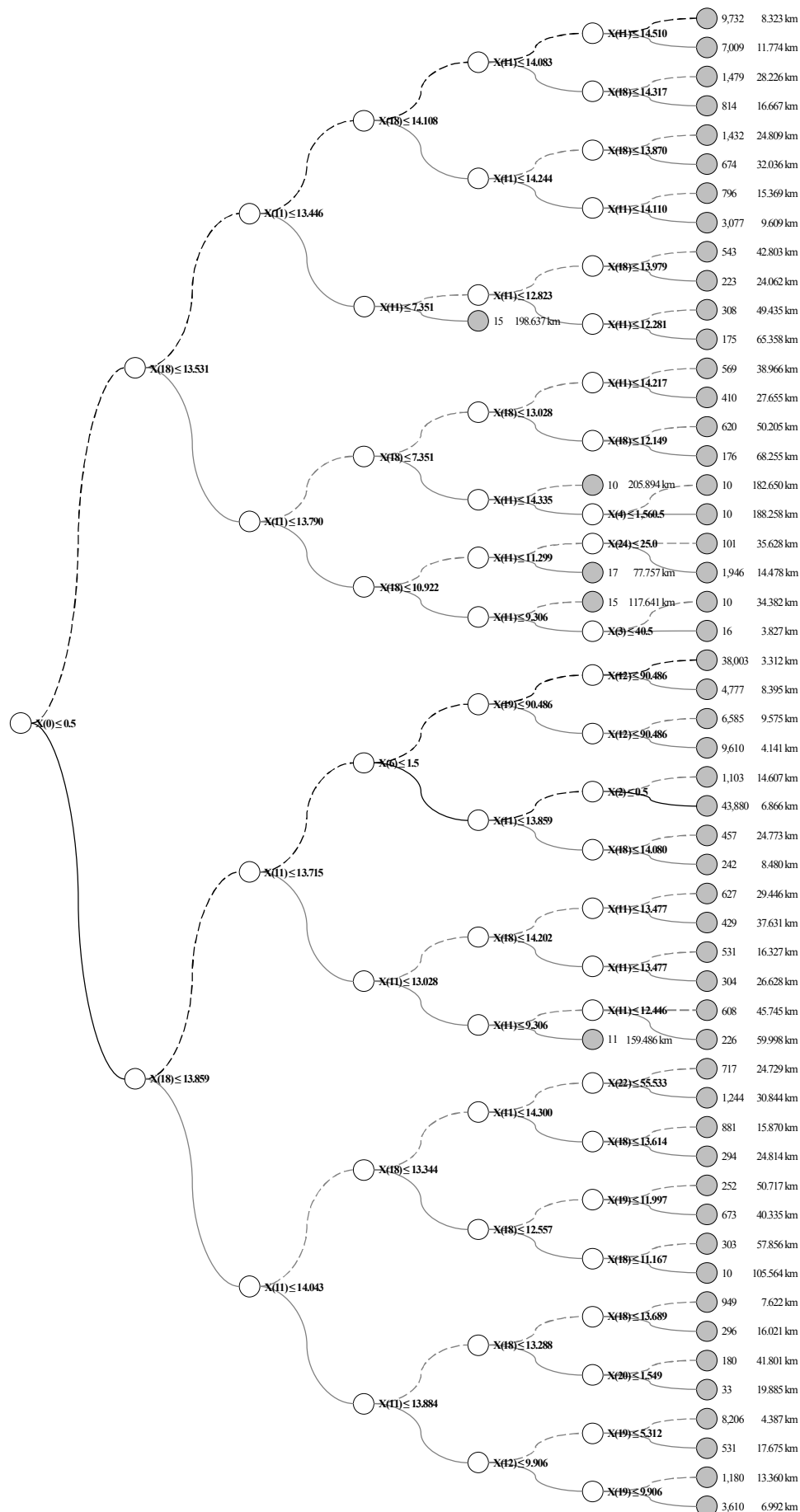
**Figure 3.** Decision Tree for Residential Relocation Distance of the Households in SMR

409

410    *5.2. Application of Ordinary Least Squares Regression and Decision Tree Regression Models*

411    This study focuses on not only the identification of the features and their structures affecting
412    residential relocation distance but also on the applicability of the machine learning approach to
413    residential mobile pattern analysis. Therefore, the application results of the previously constructed
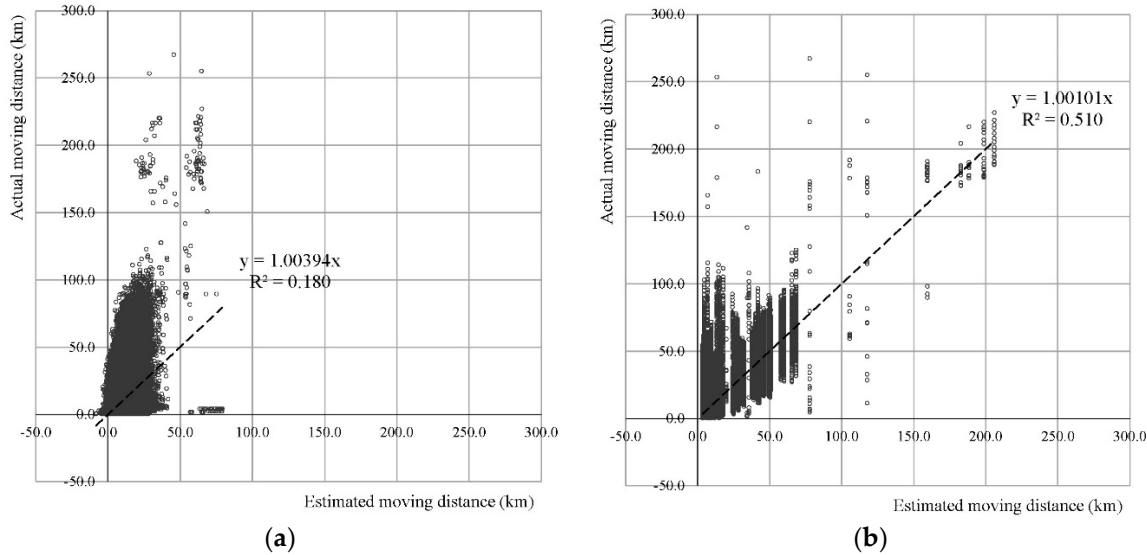414    regression models and the actual moving distance values were directly compared.
415



(**a**)                                                                (**b**)

416    **Figure 4.** (**a**) Result of Applying the Ordinary Least Squares Regression Model; (**b**) Result of
417                                  Applying the Decision Tree Regression Model

418    Figure 4 shows the comparison between the application results. The figure on the left is the result
419    of comparing the actual moving distances of all samples and the estimated moving distance using
420    the ordinary least squares regression model. The figure on the right is the comparison of the actual
421    distance and the estimated distance by the decision tree regression model. As expected, the decision
422    tree regression model results were relatively better. The application of the ordinary least squares
423    regression model showed a large number of underestimated values and a large number of unrealistic
424    residential moving distances, such as values less than zero. Whereas, the results of the decision tree
425    application presented relatively few underestimated values, and there were no unrealistic estimates
426    of the residential relocation distance. Thus, in the latter, the regression coefficient (1.00101) and the
427    R-squared value (0.510) were also better.

428    **6. Summary and Concluding Remarks**

429    In the rapidly changing Korean housing market, from both supply and demand perspectives,
430    understanding the spatial patterns of residential relocation is a meaningful task. This paper focused
431    on the structure among determinants affecting residential relocation distance and the applicability of
432    a new approach using spatial Big Data and a machine learning methodology. The results of the
433    empirical analysis on residential relocation distance in the SMR by using ordinary least squares and
434    decision tree regressions can be summarised as follows.
435    In terms of explanatory power, the decision tree regression model showed better performance
436    than the ordinary least squares regression model. Twenty variables were significant in the ordinary
437    least squares regression, whereas only twelve features were applied in the decision tree regression
438    model, although the model had relatively complicated structures. As a result of the ordinary least
439    squares regression, residential movements for housing-related reasons were shorter than the distance
440    of residential movements caused by job or education. Households with a householder over 60-years-
441    old or male householders showed longer residential relocation distance. On the other hand,
442    households with a householder less than 60-years-old, households with multiple members, and
443    households with school-aged children moved to a relatively close residential districts. In terms of the

444 location characteristics in the origin and destination, accessibility to employment markets and
445 housing ownership were the factors that shortened the household residential relocation distance. In
446 the origin, the high population density led to longer residential movements, and the variables
447 associated with the proportion of new buildings were factors that shortened the residential moving
448 distance. However, those in the destination had the opposite effects.

449      To summarise the main outcomes of the decision tree regression, the most important features
450 that determined the residential relocation distance were migration caused by a job and accessibility
451 to employment markets, although a large number of residential relocations occurred for reasons other
452 than a job. Additionally, this empirical study showed many residential movements to the districts
453 with good access to employment. The shortest moving distance was found when the household with
454 more than two people moved among densely-populated districts, whereas residential movements
455 caused by job had a relatively longer moving distance.

456      Moreover, the ordinary least squares regression and the decision tree regression models were
457 applied to compare their estimated values and the actual measurements based on the geographic
458 data using the microdata of the Internal Migration Statistics. The estimated distances using the
459 decision tree regression model were more realistic, with the estimated moving distances not
460 containing values less than zero and there were few underestimated values. Its explanatory power
461 was higher than that of the ordinary least squares regression model.

462      Thus, this study reviewed the applicability of the machine learning method using spatial Big
463 Data, which is a focus in the urban planning and management field. In particular, this article
464 attempted to overcome the limitation of conventional statistical models—low explanatory power and
465 a lot of rigid constraints—using an interpretable and understandable machine learning model, the
466 decision tree regression model. The results of this study have the following implications. First, the
467 result of the decision tree regression model (the training R-squared: 0.512) showed a significant
468 improvement in the explanatory power compared to that of the ordinary least squares regression
469 model (the training R-squared: 0.180), which is similar to a conventional linear regression model.
470 Second, the derived decision tree presented not only the diversity of structures that determine the
471 residential relocation distance, but also the main features, such as movement caused by jobs and
472 accessibility to employment markets, which form the structures. Finally, for the residential moving
473 pattern, we found that the machine learning approach, such as decision trees, can estimate more
474 realistic results than conventional methodologies.

475      The development of the forecasting model beyond the empirical analysis of the decision
476 structures for the residential relocation distance and the inclusion of several explanatory variables
477 that were not contained in the model require further research. In spite of these future tasks, this study
478 presents a case using the machine learning approach with spatial bigdata in the urban planning and
479 management field. Moreover, the outcomes of this research provide significant information about the
480 sustainable urban management of metropolitan residential districts and the construction of
481 reasonable housing policies, and it is expected to be the basis of further studies on spatial patterns of
482 residential relocation in the future.

485 **Appendix A**

486      Regarding the selection of the appropriate parameters of the decision tree for controlling the
487 overfitting problem, a trial and error method was applied in this study. In Figure A, the figure on the
488 left represents the variation of the mean squared error value (MSE) according to the leaf node
489 minimum sample value. The variation in the MSEs according to the leaf node minimum sample value
490 without setting the depth of the tree shows the lowest level from 9 to 12 samples. The figure on the
491 right indicates the variation of the R-squared values according to depth of tree after the leaf node
492 minimum sample value is set to 10. The R-squared values, which mean the explanatory powers of
493 the models applied to the training dataset and the test dataset, show the gaps of less than 1% up to a

494     depth of 6. Therefore, the appropriate parameter values of this decision tree regression model were
495     selected as follows: the leaf node's minimum sample value is 10 and the maximum depth is 6.
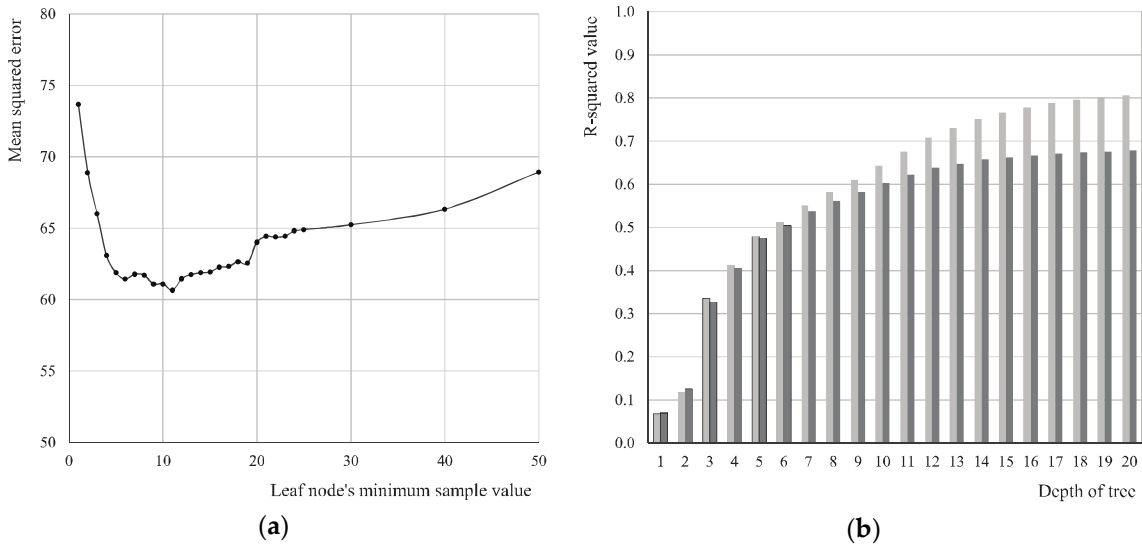


(**a**)                                              (**b**)

496     **Figure A.** (**a**) Variation of the mean squared errors according to the leaf node minimum sample
497     value; (**b**) Variation of the R-squared values according to depth of tree

## References

499     1.     Seo, D., Kwon, Y. In-migration and housing choice in Ho Chi Minh City: Toward sustainable housing
500            development in Vietnam. *Sustainability* **2017**, *9*, 1–17, doi:10.3390/su9101738.
501     2.     Rossi, P. H. *Why families move: A study in the social psychology of urban residential mobility*; Free Press:
502            New York, 1955;
503     3.     Hoyt, H. *The Structure and Growth of Residential Neighborhoods in American Cities*; Washington, 1939;
504     4.     Ravenstein, E. G. The Laws of Migration. *Journal of the Statistical Society of London* **1885**, *48*, 167–235.
505     5.     Park, E. .; Lee, J. W. A study on policy literacy and public attitudes toward government innovation-
506            focusing on Government 3.0 in South Korea. *Journal of Open Innovation: Technology, Market, and
507            Complexity* **2015**, *1*, 1–13, doi:10.1186/s40852-015-0027-3.
508     6.     Nho, H. J. Research ethics education in Korea for overcoming culture and value system differences.
509            *Journal of Open Innovation: Technology, Market, and Complexity* **2016**, *2*, 1–11, doi:10.1186/s40852-016-0030-
510            3.
511     7.     Yi, C.; Lee, S. An empirical analysis of the characteristics of residential location choice in the rapidly
512            changing Korean housing market. *Cities* **2014**, *39*, 156–163, doi:10.1016/j.cities.2014.03.002.
513     8.     OECD *OECD Regions at a Glance 2016*; 2016; ISBN 9789264252097.
514     9.     Kwon, W. Y. An examination of residential location behavior in the Seoul metropolitan area. *The
515            Annals of Regional Science* **1984**, *18*, 33–48.
516     10.    The Seoul Institute. *Seoul Statistical Series Section 04. Housing*; Seoul, 2014;
517     11.    Hong, S., Lee, Y. Limitation of Residential Mobility Distance in Seoul Metropolitan Area -Focused on
518            Migration Region and Family Size-. *Korea Real Estate Academy Review* **2015**, *60*, 115–126.
519     12.    Brown, L. A.; Moore, E. G. The inter-urban migration process: A perspective. *Geografiska Annaler Series
520            B, Human Geography* **1970**, *52*, 1–13.
521     13.    Chun, H. S. The characteristics of housing mobility of the residents' in new town areas. *Gyeonggi Forum*
522            **2004**, *6*, 91–111.

523     14.     Wu, F. Intraurban residential relocation in Shanghai: Modes and stratification. *Environment and*
524             *Planning A* **2004**, *36*, 7–25, doi:10.1068/a35177.

525     15.     Moore, E. G. *Residential mobility in the city*; 1970;

526     16.     Loren, D. L.; Eva, K.; Boaz, K. Residential relocation of amenity migrants to Florida: "Unpacking" post-
527             amenity moves. *Journal of Aging and Health* **2010**, *22*, 1001–1028, doi:10.1177/0898264310374507.

528     17.     Lee, B. H. Y.; Paul, W. Residential mobility and location choice: A nested logit model with sampling of
529             alternatives. *Transportation* **2010**, *37*, 587–601, doi:10.1007/s11116-010-9270-4.

530     18.     Stapleton, C. M. Reformulation of the family life-cycle concept: implications for residential mobility.
531             *Environment and Planning A* **1980**, *12*, 1103–1118, doi:10.1068/a121103.

532     19.     Pickvance, C. G. Life Cycle, Housing Tenure and Residential Mobility: A Path Analytic Approach.
533             *Urban Studies* **1974**, *11*, 171–188, doi:10.1080/00420987420080331.

534     20.     Morris, E. W.; Crull, S. R.; Winter, M. Housing Norms , Housing Satisfaction and the Propensity to
535             Move. *Journal of Marriage and Family* **1976**, *38*, 309–320.

536     21.     Chevan, A. Family growth, household density, and moving. *Demography* **1971**, *8*, 451–458,
537             doi:10.2307/2060682.

538     22.     Eluru, N.; Sener, I. N.; Bhat, C. R.; Pendyala, R. M.; Axhausen, K. W. Understanding Residential
539             Mobility: Joint Model of the Reason for Residential Relocation and Stay Duration. *Transportation*
540             *Research Record* **2009**, *2133*, 64–74, doi:10.3141/2133-07.

541     23.     Zanganeh, Y.; Hamidian, A.; Karimi, H. The Analysis of Factors Affecting the Residential Mobility of
542             Afghan Immigrants Residing in Mashhad ( Case Study : Municipality Regions 4 , 5 and 6 ). *Asian Social*
543             *Science* **2016**, *12*, 61–69, doi:10.5539/ass.v12n6p61.

544     24.     Ha, S. K. *Housing policy and practice in Korea*; 3rd ed.; Pakyoungsa: Seoul, 2006;

545     25.     Min, B.; Byun, M. Residential Mobility of the Population of Seoul: Spatial Analysis and the
546             Classification of Residential Mobiltiy. *Seoul Studies* **2017**, *18*, 850102.

547     26.     Clark, W. A. V.; Onaka, J. L. Life Cycle and Housing Adjustment as Explanations of Residential
548             Mobility. *Urban Studies* **1983**, *20*, 47–57.

549     27.     Yi, C.; Lee, S. I. Analyzing the Factors on Residential Mobility According to the Household Member's
550             Change - In consideration of residential duration of the households in the Seoul Metropolitan Area.
551             *Journal of Korea Planning Association* **2012**, *47*, 205–217.

552     28.     Yee, W.; Van Arsdol, M. D. J. Residential Mobility , Age , and the Life Cycle 1. *Journal of Gerontology*
553             **1977**, *32*, 211–221.

554     29.     Clark, W. A. V. Life course events and residential change : unpacking age effects on the probability of
555             moving. *Journal of Population Research* **2013**, *30*, 319–334, doi:10.1007/s12546-013-9116-y.

556     30.     Burnley, I. H.; Murphy, P. A.; Jenner, A. Selecting Suburbia : Residential Relocation to Outer Sydney.
557             *Urban Studies* **1997**, *34*, 1109–1127.

558     31.     Yang, J. Transportation Implications of Land Development in a Transitional Economy: Evidence from
559             Housing Relocation in Beijing. *Transportation Research Record* **1996**, *1954*, 7–14.

560     32.     Clark, W. A. V.; Dieleman, F. M. *Households and housing: choice and outcomes in the housing market*; New
561             Brunswick: New Jersey, 1996;

562     33.     Dieleman, F. M. Modelling residential mobility ; a review of recent trends in research. *Journal of*
563             *Housing and the Built Environment* **2001**, *16*, 249–265.

564     34.     Shalev-Shwartz, S.; Ben-david, S. *Understanding Machine Learning : From Theory to Algorithms*; 1st ed.;
565             Cambridge University Press: New York, 2014; ISBN 9781107057135.

566    35.    Witten, I. H.; Frank, E. *Data Mining: Practical Machine Learning Tools and Techniques*; 2nd ed.; Morgan
567            Kaufmann: San Francisco, 2005; ISBN 0120884070.

568    36.    Zhang, H.; Wu, P.; Yin, A.; Yang, X.; Zhang, M.; Gao, C. Prediction of soil organic carbon in an
569            intensively managed reclamation zone of eastern China: A comparison of multiple linear regressions
570            and the random forest model. *Science of the Total Environment* **2017**, *592*, 704–713,
571            doi:10.1016/j.scitotenv.2017.02.146.

572    37.    Estelles-lopez, L.; Ropodi, A.; Pavlidis, D.; Fotopoulou, J.; Gkousari, C.; Peyrodie, A.; Panagou, E.;
573            Nychas, G. J.; Mohareb, F. An automated ranking platform for machine learning regression models for
574            meat spoilage prediction using multi-spectral imaging and metabolic profiling. *Food Research*
575            *International* **2017**, *99*, 206–215, doi:10.1016/j.foodres.2017.05.013.

576    38.    Chagas, C.; Junior, W.; Bhering, S.; Filho, B. Spatial prediction of soil surface texture in a semiarid
577            region using random forest and multiple linear regressions. *Catena* **2016**, *139*, 232–240,
578            doi:10.1016/j.catena.2016.01.001.

579    39.    Oliveira, M.; Gama, J. An overview of social network analysis. *WIREs: Data Mining and Knowledge*
580            *Discovery* **2012**, *2*, 99–105.

581    40.    Tso, G. K. F.; Yau, K. K. W. Predicting electricity energy consumption: A comparison of regression
582            analysis , decision tree and neural networks. *Energy* **2007**, *32*, 1761–1768,
583            doi:10.1016/j.energy.2006.11.010.

584    41.    Xu, M.; Watanachaturaporn, P.; Varshney, P. K.; Arora, M. K. Decision tree regression for soft
585            classification of remote sensing data. *Remote Sensing of Environment* **2005**, *97*, 322–336,
586            doi:10.1016/j.rse.2005.05.008.

587    42.    Prasad, A. M.; Iverson, L. R.; Liaw, A. Newer Classification and Regression Tree Techniques: Bagging
588            and Random Forests for Ecological Prediction. *Ecosystems* **2006**, *9*, 181–199, doi:10.1007/s10021-005-
589            0054-1.

590    43.    Hansen, W. G. How accessibility shapes land use. *Journal of the American Institute of Planners* **1959**, *25*,
591            73–76.

592    44.    Wilson, A. G. *Entropy in Urban and Regional Modeling*; Pion: London, 1970;

593    45.    Yi, Y.; Kim, E.; Choi, E. Linkage among school performance, housing prices, and residential mobility.
594            *Sustainability (Switzerland)* **2017**, *9*, 1–18, doi:10.3390/su9061075.

595    46.    Choi, Y.; Yim, H. K. Determinants of the Residents' Settlements Employing Poission Regression. *The*
596            *Korea Spatial Planning Review* **2005**, *46*, 99–114.

597    47.    Pagliara, F.; Simmonds, D. Conclusions. In *Residential location choice-models and applications*; Pagliara, F.,
598            Preston, J., Simmonds, D., Eds.; Springer, Verlag: Berlin, Heidelberg, 2010; pp. 243–248.

599    48.    Simmonds, D. The DELTA residential location model. In *Residential location choice-models and*
600            *applications2*; Pagliara, F., Preston, J., Simmonds, D., Eds.; Springer, Verlag: Berlin, Heidelberg, 2010;
601            pp. 77–97.

602    49.    Waddell, P. Modelling residential location in UrbanSim. In *Residential location choice-models and*
603            *applications*; Pagliara, F., Preston, J., Simmonds, D., Eds.; Springer, Verlag: Berlin, Heidelberg, 2010; pp.
604            165–180.

605    50.    Cadwallader, M. *Migration and Residential Mobility: Macro and Micro Approaches*; University of
606            Wisconsin Press: Madision, 1992;

607    51.    Cluttons *Residential Mobility in London: Unlocking Migration Patterns*; 2017;

608    52.    Mueller, A. C.; Guido, S. *Introduction to Machine learning with python*; O'Reilly: Sebastopol, 2016; ISBN
609           9781491917213.
610