# 戸建住宅価格における機械学習を用いた2段階推計モデル

Two-Step estimation model using machine learning at detached house price

高橋 佑典 \*1 Yusuke Takahashi

\*<sup>1</sup>富士通クラウドテクノロジーズ株式会社 FUJITSU CLOUD TECHNOLOGIES LIMITED

This paper reports the result of trying to estimate the detached house price with a two step model. By handling the advantages of linear model and nonlinear model in combination, we can expect explanation possibility and estimation accuracy for model when estimating house price. Experiments compare the methods by machine learning and the effects of explanatory variables to be input to the nonlinear model.

# 1. はじめに

住宅価格の分析を行った従来の研究では、多くが線形モデル として推定している.その背景として、シンプルなモデルであ ればモデルの推計価格に対して明確な説明が可能であるという 点が大きい.しかし、専有面積や築年数などでは、住宅価格に 与える影響が区間によって異なり、非線形になると考えられ、 従来の研究でもその存在が示されている.また、線形モデルと して推定する場合、住宅価格を決定する上で主要な要因である 築年数と建築年代といった多重共線性を持つ説明変数を同時に 投入してしまうと、モデルの推計値が不安定になってしまうこ とも挙げられる.

そこで, 先行研究のうち [Shimizu 14] では, 非線形モデルと して推計するため, ノンパラメトリックなモデルである連続量 dummy モデル (DmM) と, AIC を評価指標とした Switching RegressionModel (SWR), 一般加法モデル (GAM) が用い られ, 非線形性を考慮した推計が行われている. 非線形モデル を扱う課題として, モデルが複雑になり再現性が低下すること が考えられる.

本研究では、線形モデルと非線形モデルの課題に対処するた め、モデルを2段階に分けて戸建住宅における取引価格の推 計を行う.1段階目のモデルでは線形モデルとして推計し、モ デルに対する説明性を担保する.2段階目のモデルでは1段階 目の誤差項を目的変数として、非線形なモデルによって推計を 行う.本研究の最終的な目的は、線形モデルと非線形モデルを 組み合わせた二段階のモデルによって推計することで、住宅価 格の推計時に、説明性と推計精度を向上させることである.

実験では、最初に線形モデルと機械学習による非線形モデルを作成し、誤差率の分布を比較することで、非線形性をうまく表現できる手法の選択を行う.次に、研究背景で述べた非線形性の考慮を機械学習を用いて行い、先行研究で扱われている一般加法モデル(GAM)の結果と比較する.評価指標としては機械学習で用いられている指標である平均平方二乗誤差(RMSE)、平均平方二乗誤差率(RMSPE)を用いる.

## 2. 関連研究

## 2.1 非線形性

[Shimizu 14] では東京 23 区の中古マンションを対象に,線 形モデルをベースに置き,物件の平米単価と各説明変数間の関 係を非線形モデルとして推計している.具体的には,ノンパラ メトリックなモデルである連続量 dummy モデル (DmM)と, AIC を評価指標とした Switching Regression Model (SWR), 一般加法モデル (GAM)を用いて中古マンションの取引価格 の主要要因が持つ非線形性を明らかにしている.

[小野 02] ではリクルート社の「週刊住宅情報」に掲載され た東京 23 区の中古マンションデータを扱い,中古マンション 価格に影響を与える要因のうち,「建築後年数」と「最寄り駅 までの徒歩時間」に非線形性が存在することを明らかにして いる.

#### 2.2 機械学習

[福井 18] では中古マンションの取引データをもつ,レイン ズの成約データに対してニューラルネットワークを用いて,不 動産査定価格を分類問題として解いている.

[大和 18] では不動産情報サイト SUUMO におけるデータを 扱い,家賃の推計を行う際に変数の非線形性を表現できる手法 として,線形回帰モデルの代わりに決定木をベースとしたラン ダムフォレストを用いている.

本研究では、関連研究で明らかにされている、住宅価格に影響を与える要因に非線形性が存在する点と、それらの非線形性 を機械学習によって表現し、推計価格の精度向上が図れるかを 調査する.また、先行研究では主に中古マンションの取引価格 データが扱われていることに対し、本研究では戸建住宅の中古 取引価格データを扱う.

## 3. 実験

## 3.1 データセット

本研究においては、国土交通省が提供している 土地総合情報システムのうち、不動産取引価格情 報ダウンロードサービス([土地総合情報システム](http://www.land.mlit.go.jp/webland/download.html)) から取得したデータを用いた.取得したデータは、取引時期 が2005年第3四半期から2018年第3四半期である東京23 区内の取引価格情報である.本研究では、戸建住宅の中古取 引価格データを扱いたいため、データの種類が「宅地(土地

連絡先: 高橋 佑典, 富士通クラウドテクノロジーズ株式会社,
東京都中央区銀座7丁目16番12号 G-7ビルディング,
03-6281-5710, yus-takahashi@fujitsu.com

と建物)」,用途が「住宅」,地域が「住宅地」となっている レコードのみを抽出した.この不動産取引情報には,不動産 取引において一般的な取引時期,延床面積,土地面積,最寄 り駅までの徒歩分数といった情報が取引価格とともに含まれ る.表1に本研究で扱うデータの要約統計量を示す.各デー タはそれぞれ TS:最寄り駅までの徒歩分数,P:成約価格,L: 土地面積,S:延床面積,RW:前面道路幅員,BLR:建ペい率, FRA:容積率,CY:建築年,A:建築後年数,WOOD:木造ダ ミーである.また,要約統計量で示したデータ以外にも用途 地域や市区町村コードといったデータも含まれている.

- <u>F</u> -			2	~	コピカレカセント	Ħ
÷	1.	エー	1	(1)	男約除計	<b>音</b>
1	±	/	<u></u>	~ /	シャコルルロト	440.

	count	mean	std	min	25%	50%	75%	max
TS	14424	10.65654	5.411991	1	7	10	14	29
P	14424	48846850	17609200	21000000	35000000	46000000	6000000	99000000
L	14424	93.0633	31.69272	50	70	90	110	200
s	14424	97.72463	24.75509	50	85	95	105	200
RW	14424	4.693795	1.288334	2.5	4	4	5.5	9
BLR	14424	56.56129	7.888145	30	50	60	60	80
FRA	14424	174.0294	76.58877	60	100	150	200	800
CY	14424	1996.819	13.50506	1955	1988	2000	2007	2017
TDY	14424	2011.655	3.736598	2005	2008	2012	2015	2018
A	14424	14.83555	13.07827	1	2	12	23	50
WOOD	14424	0.8840821	0.3201376	0	1	1	1	1

#### 3.2 二段階推計モデル

本研究では、モデルを二段階に分けて推計している.理由と して、1段階目に線形モデルを推計して説明性を担保しつつ、 2段階目のモデルでは多重共線性を持ってしまう要因や非線形 性を持つと思われる要因を投入し、1段階目の誤差を予測し補 正することで、モデル全体としての予測精度の向上を図る. step1.

 $log(y) = \alpha + \beta_1 a_1 + \beta_2 a_2 + \beta_3 a_3 + \dots + \mu_1$ step 2.

 $\mu_{1} = \gamma + \delta_{1}b_{1} + \delta_{2}b_{2} + \delta_{3}b_{3} + \dots + \mu_{2}$ 

線形モデルを二段階に分けて推計する場合,上記のような2 つの線形モデルが作られる.本研究では,1段階目では線形モ デルを推計し取引価格を予測する.2段階目では1段階目の予 測誤差を目的変数として,非線形なモデルを推計し,取引価格 に対する最終的な予測を行う.

本研究では、戸建住宅の取引データに対して一般加法モデ ル(GAM)を用いて住宅価格に影響を与える主要な要因の非 線形性を確認した.一般加法モデル(GAM)を適用した結果 を図1に示す.

ー般加法モデル(GAM)を適用した結果から,極端な傾き の変化は見られなかったものの,S:延床面積,L:土地面積,A: 建築後年数,といった主要な要因に対して,非線形性を確認す ることができた.

#### **3.3** 手法の比較

本研究では,線形モデルである OLS をベースラインとし, 比較対象として,機械学習の手法である決定木をもとにした ランダムフォレストと xgboost を選択した.また,ニューラ ルネットワークによる手法として多層パーセプトロンを選択し た.以下で各手法についての概要を説明する.

OLS は定数項(切片)と説明変数の係数によって値を予測 する線形モデルであり,最小二乗法によって係数と切片を決定 する.ランダムフォレストは,決定木とバギングを組み合わせ た手法であり,決定木を大量に生成し、各決定木の結果を集計



図 1: 各変数と取引価格の関係:GAM

して予測を行う.各決定木は独立して異なる特性を持つように 学習する.xgboost は,決定木とブースティングを組み合わせ た手法であり,決定木を逐次的に増やしていき、生成済みの決 定木の誤差を補正するように、新たな決定木を生成し学習を進 めていく.多層パーセプトロンは,入力、中間、出力の3層か らなるニューラルネットワークの手法。バックプロパゲーショ ンを用いた学習を行う.

機械学習においては、学習に使ったデータセットだけに過度 に適合したパラメータが学習されてしまい、テストデータに 対して性能が出ない過学習と呼ばれる状態に陥ることがある。 過学習を抑えるために、決定木による手法では特徴量の数や生 成する木の深さを設定することが考えられる.また NN によ る手法では、学習の過程において大きな重みを持つことに対し てペナルティを課す Weight decay と呼ばれる正則化の手法が 存在する.

本研究における機械学習手法では、グリッドサーチによる パラメータチューニングを行い、過学習を抑制するためにパラ メータを定めてからモデルの構築を行った.

#### 3.4 評価

まず,手法間の誤差率分布を比較した結果を図2に示す.デー タは学習データとテストデータを8:2でランダムに分割し,各 手法において同じ説明変数を投入したモデルで一つ一つの物件 に対して取引価格の予測を行い,予測価格/実際の取引価格で 定義される誤差率を算出した.その実験を200回繰り返し行 い,誤差率の平均値の分布を作成した.

結果から,OLS や決定木をベースにした手法は予測結果が 大きくはずれてしまう可能性が少ないことに対して,ニューラ ルネットワークをベースにした多層パーセプトロンでは誤差率 の平均値の分布に外れ値が見られた.このことから多層パーセ プトロンでは推計結果に大外れが生じてしまうことが確認で きた.

#### **3.5** 2 段階目投入変数

第3.4節の結果から,2段階目の手法として大外れが少な く,過学習も抑制できていると考えられた xgboost を選択し た.2段階目の機械学習モデルに投入する変数として,他の 説明変数への影響が明確でない STATION:最寄り駅ダミーや CENSUS:大字・通称レベルの地域ダミーを投入してその効果 を確認した.また,S:延床面積,L:土地面積,A:建築後年数, TS:最寄り駅までの徒歩分数を説明変数として単独投入し,そ の効果を確認した.結果を箱ひげ図として図3と図4.に示す. また,各評価指標の値としては表4.と表3に示す.

これらの結果から評価指標を見ると,機械学習における非 線形性の考慮が,僅かだが精度向上に寄与していることが確認 できた.

## 4. おわりに

本研究では、戸建住宅データにおける取引価格に影響を与え る主要要因に対して、機械学習での非線形性を考慮した2段 階推計を行った.実験結果より、手法間の比較では、ニューラ ルネットワークをベースにした手法よりも、決定木をベースに した手法の方が、誤差率のばらつきが少ないことがわかった. また、2段階推計においては、2段階目のモデルに個別に説明 変数を投入し、その補正結果を調べた.結果は一般加法モデ ル(GAM)で確認した非線形性を考慮し、精度向上が確認で きた.また、合わせて戸建住宅価格に影響を与える最寄り駅ダ ミーと大字・通称ダミーを二段階目の説明変数として投入し、 精度向上を確認できた.今後は、今回得られた手法ごとの誤差 率分布をもとに、非線形なモデルを推計する手法ごとの時性を 明らかにすることや、モデルに投入する説明変数の数を増やし ていき、その効果を検証することが考えられる.



図 2: 誤差率の平均値の分布



図 3: ダミー変数の効果



図 4: 説明変数ごとの効果

表 2:2 段階目投入変数別の効果(ダミー変数)

	RMSE	RMSPE
BaseModel	10038752	5.1895
STATION	9566984	4.7900
CENSUS	9364350	4.6412

表 3:2 段階目投入変数別の効果(連続量)

	RMSE	RMSPE
BaseModel	10038752	5.1895
S	9942438	5.1025
L	9976350	5.1309
А	9902803	5.1253
TS	10033856	5.1733

## 参考文献

- [Shimizu 14] Shimizu, C., Karato, K. and Nishimura, K.: Nonlinearity of housing price structure assessment of three approaches to nonlinearity in the previously owned condominium market of Tokyo, Int. J. Housing Markets and Analysis, Vol. 7, No. 4, pp. 459-488 (2014).
- [小野 02] 小野宏哉,高辻秀興,清水千弘:「品質を考慮した中 古マンション価格モデルの推定」,麗澤経済研究,第10 巻第2号, pp.81-102, 2002.
- [福井 18] 福井光,阪井一仁,南村忠敬,三尾順一,木下明弘, 田司郎:「レインズのニューラルネットワークを用いた不動 産価格査定について」, The 32nd Annual Conference of the Japanese Society for Artificial Intelligence, 2018.
- [大和 18] 大和大祐,野村眞平:「SUUMO でのビッグデータ活 用事例」,日本不動産学会誌/第 31 巻第 1 号, pp.78-83, 2017.