

# データを一から分析してみよう(Part 2)

Kaggle で公開されているアイオワ州住宅価格のデータセットを使って、EDA (探索的デー タ解析)からモデル作成までのデータ分析プロセスを行なっていきます。Part1 は、データ の読み込みと正しいデータ型へ修正するところまで行いました。Part2 では、データを確認 しながら欠損値や外れ値の処理を行い、モデルを作成していきます。Kaggle のコンペティ ションに実際に参加し、作成したモデルの結果 (順位)を確認するまでを掲載しています。 ※本資料は、macOS、RapidMiner Studio Version9.10 を使用しています。



使用データ: Ames Housing dataset
 アメリカ合衆国アイオワ州エイムズ都市の住宅に関するデータセット
 <a href="https://www.kaggle.com/c/house-prices-advanced-regression-techniques/overview">https://www.kaggle.com/c/house-prices-advanced-regression-techniques/overview</a>
 外装素材、地下室の高さ、ガレージのサイズなど、79の属性(説明変数)
 各属性(説明変数)の特徴は data\_description.txt を参照

分析の目的: 各住宅の販売価格を予測すること

### <Part2 で行うプロセス>

- データの理解
- 欠損値の処理
- 目的変数との相関
- 説明変数同士の相関
- 外れ値の処理
- モデル作成・適用
- Kaggle のコンペティションへ参加



データの理解

Part1 で作成した以下のプロセスを実行し、データの特徴や欠損値の有無などを確認します。



▽ 目的変数(SalePrice)のヒストグラムから、データの特徴を把握します



データの約 60%が 100,000 ドル~200,000 ドルを占めています。

v						
	属性名	・・ データ型	欠損値	基本統計量		
~	LotFrontage	Numeric	259	最小值 21	最大值 313	平均值 70.050
	🔥 GarageYrBlt	Numeric	81	最小值 1900	最大值 2010	平均值 1978.506
~	MasVnrArea	Integer	8	最小值 O	最大值 1600	平均值 103.685

#### ▽ 欠損値の有無を確認します



3つの属性に欠損値がありました。次のステップで欠損値の処理を行います。 LotFrontage:敷地に接続する道路の直線距離 GarageYrBlt:ガレージが建てられた年 MasVnrArea:石積みベニヤの面積(平方フィート)

欠損値の処理			
Process Retrieve ames-housi Numerical to Polyno Parse Numbers Dipp out exa exa ori exa ori ori exa ori	1 Filter Examples exa ア exa ori umm LotFrontage属性 欠損値を含まない 行のみ選択	2 Replace Missing Val ( exa  exa  ori pre 欠損値を0に置換	res res

Step1. LotFrontage 属性は、欠損値が多いので、Filter Examples オペレータを使用して 欠損値を含まない行のみを出力するように設定します。

<パラメータの設定>





Step2. GarageYrBlt 属性は、GarageType 属性(ガレージの場所)を確認すると None(ガレージなし)が 74 あり、ガレージ自体が存在しないため欠損していたことがわかりました。MasVnrArea 属性も、MasVnrType 属性(ベニヤタイプ)を確認すると None(ベニヤなし)が 6 あり、ベニア自体が存在しませんでした。

→Replace Missing Values オペレータを使用し、欠損値を 0 に置き換えます。

パラメータ ×	環境XX
👫 Replace Missing Val	ues
create view	٦
attribute filter type	no_missing_values 🔻 🗊
✓ invert selection	٩
include special attrib	outes (1)
default	zero 🔻 🗊

欠損している値を zero に置き換えたいので、invert selection にチェックを入れます。 設定後プロセスを実行すると、欠損値が全て 0 に変換されていることが確認できます。

### 目的変数との相関

属性数が 79 と多いので、目的変数と各属性の相関係数を求め、説明変数の重要度を確認し ます。



- Step0. 欠損値処理までのプロセスをサブプロセスにまとめました。
- Step1. 相関係数を計算するには、全て数値型に変更してから行う必要があります。 Nominal to Numerical オペレータ使用し、カテゴリ型から数値型へ変更します。



<パラメータの設定>

パラメータ ×	環境 ×
ᅌ Nominal to Numeric	al
create view	٦
attribute filter type	all 🔹 🗊
invert selection	٢
include special attril	outes
coding type	unique integers 🔻 🗊

coding type は unique integers を指定することで、カテゴリ属性の値が等しくランク付けさ れた実数値の属性に変更されます。

Step2. Weight by Correlation オペレータを使用し、目的変数との相関係数から説明変数 の重要度を確認します。

	結果概要	📑 Attribu	teWeights (Weight by Correlation)
	データ	attribute w Overall 0	reight .802
<ハフメータの設定> パラメータ × 環境 ×	Weight	GrLivArea 0 Garage 0	.704 .647
Weight by Correlation	Visualizations	Garage 0	.632
normalize weights		TotalBs 0 1stFlrSF 0	.627
✓ sort weights	アノテーション	FullBath 0	.567
sort direction descending		YearBuilt 0	.539
		YearRe 0	.519
squared correlation		MasVnr 0	.492

sort direction は descending を指定し、降順にソートされるように設定します。 プロセスを実行すると、目的変数と各属性との相関係数を降順で確認することができます。

目的変数と相関が高い属性を説明変数とします。今回、相関係数が 0.5 以上であった 10 の 属性を選択して分析を進めていきます。(変数の選択は次のステップで行います) 説明変数:OverallQual、GrLivArea、GarageCars、GarageArea、TotalBsmtSF 1stFlrSF、FullBath、YearBuilt、TotRmsAbvGrd、YearRemodAdd



ここで注意したいのが、「多重共線性」という問題です。多重共線性とは、説明変数同士に 強い相関があるときに起こる状態のことで、分析結果の信頼性が低下してしまう可能性が あります。説明変数の選択は、多重共線性を考慮して行う必要があります。

# 説明変数同士の相関

多重共線性を防ぐために、まず説明変数同士に強い相関があるものが存在するかを確認し、 強い相関が存在する場合はいずれか一方を削除します。



(Nominal to Numerical と Weight by Correlation は属性の重要度を確認するために使用したので、無効化しました)

Step1. Select Attributes オペレータを使用し、目的変数との相関係数が 0.5 以上であった 10 の属性を説明変数として選択します。

<パラメータの設定>





Step2. Correlation Matrix オペレータを使用して相関行列を生成し、説明変数同士に強い 相関があるものが存在するかを確認します。

<パラメータの設定>

パラメータ ×	環境	×
Correlation Matrix		
attribute filter type	all	•
invert selection		٢
🗸 include special attri	butes	<b>()</b>
✓ normalize weights		٢
squared correlation	1	٢

include special attributes にチェックを入れて実行すると、特別属性である label 属性と id 属性を含む相関行列が以下のように生成されます。

Attributes	SalePrice	OverallQual	YearBuilt	YearRemodAdd	TotalBsmtSF	1 stFIrSF	GrLivArea	FullBath	TotRmsAbvGrd	GarageCars	GarageArea	ld
OverallQual	0.802	1	0.588	0.561	0.570	0.521	0.608	0.564	0.445	0.613	0.584	-0.039
YearBuilt	0.539	0.588	1	0.597	0.421	0.330	0.208	0.483	0.111	0.543	0.498	-0.012
YearRemodAdd	0.519	0.561	0.597	1	0.315	0.283	0.295	0.456	0.196	0.426	0.389	-0.024
TotalBsmtSF	0.627	0.570	0.421	0.315	1	0.833	0.478	0.340	0.302	0.458	0.512	-0.019
1 stFIrSF	0.620	0.521	0.330	0.283	0.833	1	0.573	0.390	0.418	0.464	0.511	-0.005
GrLivArea	0.704	0.608	0.208	0.295	0.478	0.573	1	0.620	0.829	0.474	0.474	-0.003
FullBath	0.567	0.564	0.483	0.456	0.340	0.390	0.620	1	0.548	0.471	0.418	0.010
TotRmsAbvGrd	0.537	0.445	0.111	0.196	0.302	0.418	0.829	0.548	1	0.378	0.350	0.020
GarageCars	0.647	0.613	0.543	0.426	0.458	0.464	0.474	0.471	0.378	1	0.890	0.003
GarageArea	0.632	0.584	0.498	0.389	0.512	0.511	0.474	0.418	0.350	0.890	1	-0.012
SalePrice	1	0.802	0.539	0.519	0.627	0.620	0.704	0.567	0.537	0.647	0.632	-0.037
ld	-0.037	-0.039	-0.012	-0.024	-0.019	-0.005	-0.003	0.010	0.020	0.003	-0.012	1

今回、説明変数同士の相関が 0.8 以上の組み合わせをピックアップし、目的変数との相関 係数が低い方を削除します。SalePrice(目的変数)列は、各説明変数と目的変数との相関係数 を表すため、ドラッグ&ドロップで Attributes 列のすぐ右隣に移動させました。

説明変数同士の相関が 0.8 以上であった組み合わせは3つありました。

- ・GarageCars と GarageArea …… 0.890
- ・TotalBsmtSF と 1stFlrSF ……… 0.833
- ・GrLivArea と TotRmsAbvGrd … 0.829

目的変数との相関係数を確認した結果、GarageArea、1stFlrSF、TotRmsAbvGrd を削除します。



多重共線性を考慮し、以下の7つを説明変数とします。 説明変数:OverallQual、GrLivArea、GarageCars、TotalBsmtSF、FullBath YearBuilt、YearRemodAdd(変数の選択は次のステップで行います。)

# 外れ値の処理

モデルの作成をする前に、データの分布をみて他の値から大きく外れた値(外れ値)の有 無を確認し、外れ値があった場合は取り除きます。



Step1. Select Attributes オペレータを使用し、7つの説明変数を選択します。(パラメータ 設定は p.6 参照)各説明変数の散布図をみて外れ値の有無を確認します。 "X-Axis column"を変更すると、各説明変数のデータ分布が確認できます。







外れ値の有無を確認した結果、以下を外れ値として処理することにします。

OverallQual ……評価が 10 以上で、住宅価格が 200,000 ドルより下は削除

TotalBsmtSF ……地下室総面積が 3000 より上で、住宅価格 300,000 ドルより下は削除

FullBath …………評価が1より下で、住宅価格が250,000 ドルより上は削除

YearBuilt………1940年以前に建てられた物件で、住宅価格が 260,000 ドル以上は削除 2000年以前に建てられた物件で、住宅価格が 60,0000 ドル以上は削除

YearRemodAdd…1980 年より前のリフォーム物件で、住宅価格が 35,0000 より下は削除

Step2. Generate Attributes オペレータを使用して、Step1 で決めた外れ値の条件を指定し、 外れ値のフラグ付けした列を作成します。

<パラメータの設定>





パラメータリストを編集: function descriptions List of functions to generate.	s
attribute name	function expressions
OverallQual_flag	if(OverallQual>=10&&SalePrice<200000,"true","false")
TotalBsmtSF_flag	if(TotalBsmtSF>3000&&SalePrice<300000,"true","false")
FullBath_flag	if(FullBath<1&&SalePrice>250000,"true","false")
YearBuilt_flag	if(YearBuilt<=1940&&SalePrice>=260000,"true","false")
YearBuilt_flag2	if(YearBuilt<=2000&&SalePrice>=600000,"true","false")
YearRemodAdd_flag	if(YearRemodAdd<1980&&SalePrice>350000,"true","false")
	📜 エントリを追加 🔡 エントリを削除 💽 適用 🗙 キャンマル
● ◎ ◎ 式を/	a編集: function expressions
Info: Expression is syntactically correct.	×
Functions 校索	Inputs 校索 下
論理演算子	A Gri ivArea
if (条件, trueの戻り値, falseの戻り値)	OverallOual
% I	① # TotalBsmtSF
% &&	① # YearBuilt
% II	①     # YearRemodAdd     ■
比較演算子	※
テキスト情報	(%)     (
テキスト変換	id
数学関数	SalePrice
統計関数	🛞 🗸 label 🗸
	🧭 適用 🗶 キャンセル

"属性名+\_flag"という新しい属性(列)を作成し、数式には外れ値とする条件を入れます。 数式は直接入力または電卓マークから設定できます。適用しプロセスを実行すると、データ セットには、新しい属性(列)が以下のように追加されます。

geCars	OverallQual_flag	TotalBsmtSF_flag	FullBath_flag	YearBuilt_flag	YearBuilt_flag2	YearRemodAdd_flag
	false	false	false	false	false	false
	false	false	false	true	false	false
	false	false	false	false	false	false
		<i>c</i> .	<i>c</i> .	<i>c</i> .	<i>c</i> .	6 J

外れ値のある行には"true"、それ以外の行は"false"と表示されます。



Step3. Filter Examples オペレータを使用し、外れ値以外の列を出力します。 <パラメータの設定>

パラメータ ×	環境 ×
🝸 Filter Examples (2	2) (Filter Examples)
filters	🍸 フィルタを追加 👔
condition class	custom_filters v
🗸 invert filter	Œ

invert filter をクリックすると、フィルタを追加 から設定した条件が反転されます。

フィルタを行 Defines the	乍成: <b>filte</b> i list of filt	r <b>s</b> ers to apply.						
OverallQual_flag	•	equals	•	true			*	×
TotalBsmtSF_flag	•	equals	•	true			*	×
FullBath_flag	¥	equals	•	true			*	×
YearBuilt_flag	▼	equals	•	true			*	×
YearBuilt_flag2	•	equals	•	true			*	×
YearRemodAdd_flag	•	equals	•	true			*	×
			-					
○ すべてにマッチする	) どれ	か一つにマッチする	🗸 条件の	事前選択	エントリを追加	✓ <u>о</u> к	<b>X</b> +	ャンセル

外れ値のフラグ付けした属性(列)に、外れ値"true"が一つでもあればその行は出力しないように設定します。

Step4. Select Attributes オペレータを使用して、説明変数を再選択します。
 Step2 で外れ値のフラグ付けした属性(列)を作成して属性が増えたので、Step1 と
 同様に説明変数を7つ選択します。
 説明変数:OverallQual、GrLivArea、GarageCars、TotalBsmtSF、FullBath

YearBuilt, YearRemodAdd

ここまでのプロセスを実行すると、外れ値の存在する行を削除したことによって、データが 1,201 行から 1,187 行になります。

```
ExampleSet (1,187 行,2 特別属性,7 通常属性)
```



# モデル作成・適用

データの準備ができたので、いくつかの機械学習アルゴリズムを試して交差検証を行い、 今回扱うデータと相性が良いのはどれかを確認します。



Step1. Cross Validation オペレータを配置し、交差検証を行います。 (実行結果を揃えるため、use local random seed にチェックを入れます)

Step2. 住宅価格の予測なので、回帰の機械学習アルゴリズムを選択しモデルを作成します。 <パラメータの設定 例>

パラメータ ×	環境XX
Random Forest	
number of trees	100
criterion	least_square 🔻 🗊
criterion maximal depth	least_square ▼ ① 10

回帰を用いるので、criterion は least\_square を指定します。(Decision Tree の場合も同様) 今回、パラメータはデフォルト値のままにします。



- Step3. Apply model オペレータで、交差検証用に作成したモデルにテスト用データを適用 します。
- Step4. Performance(Regression)オペレータを使用し、検証用に作成したモデルが未知デ ータに対してどのくらいの精度があるかを、4 つの評価指標で確認します。6 つの アルゴリズムを試し、結果を比較します。

<パラメータの設定>



#### 実行結果を下記表にまとめました。(パラメータ:least\_square 指定以外はデフォルト設定)

機械学習アルゴリズム	root mean squared error	relative error	correlation	squared correlation
Random Forest(回帰)	27055.052 +/- 2395.754	11.50% +/- 1.34%	0.937 +/- 0.018	0.878 +/- 0.034
Deep Learning	28153.327 +/- 2262.050	12.64% +/- 1.26%	0.936 +/- 0.016	0.876 +/- 0.029
Neural Net	28623.182 +/- 3020.979	13.25% +/- 2.73%	0.938 +/- 0.016	0.880 +/- 0.030
Liner Regression	33036.439 +/- 4293.221	14.82% +/- 1.15%	0.906 +/- 0.015	0.821 +/- 0.027
Decision Tree(回帰)	33461.141 +/- 2613.566	14.18% +/- 1.45%	0.905 +/- 0.020	0.819 +/- 0.036
k-NN	39293.207 +/- 5267.998	15.82% +/- 1.70%	0.866 +/- 0.024	0.751 +/- 0.043

※バージョンや環境などにより、実行結果に多少誤差がある場合があります

この中では、ランダムフォレストがデータと最も相性が良いことが分かりました。



### ▽ モデル作成

交差検証の結果から、ランダムフォレストを用いてモデルを作成します。



交差検証の Trainig フェーズに配置したものと同じオペレータ(Random Forest)をプロセ ス画面に配置します。これで、予測モデルを作成できました。

#### ▽ モデル適用

生成した予測モデルを使って、実際に各住宅の販売価格を予測します。



Step1. 予測したいデータ(テストデータ)をモデルに適用させるための準備をします。 Part1 で ID 属性を定義しリポジトリに格納しておいたテストデータを配置し、訓 練データと同様の方法で正しいデータ型に修正します。

Step2. Apply Model オペレータを使用し、テストデータを予測モデルに適用させます。

プロセスを実行すると、以下のように prediction(SalePrice)の列が追加され、テストデータ に対する住宅価格の予測値が出力されます。

結果概要	🚦 Randon	n Forest Model (I	Random Forest)	🛛 🧧 Examp	leSet (Apply M	odel) ×					
	開く Tur	bo Prep	Auto Model					フィルタ (1,45	9 / 1,459 行):	all	
データ	Row No.	ld	prediction(SalePrice)	LotFrontage	GarageYrBlt	MSSubClass	MSZoning	LotArea	Street	Alley	I
	1	1461	119679.437	80	1961	20	RH	11622	Pave	NA	Ē
Σ	2	1462	149361.923	81	1958	20	RL	14267	Pave	NA	1
基本統計量	3	1463	177780.634	74	1997	60	RL	13830	Pave	NA	1
	4	1464	182618.892	78	1998	60	RL	9978	Pave	NA	ī
<b>S</b>	5	1465	210578.093	43	1992	120	RL	5005	Pave	NA	1

この予測結果を Kaggle に提出し、コンペティションに実際に参加してみます。



### Kaggle のコンペティションへ参加



Kaggle は、世界中のデータサイエンティストたちが集まり、構築したモデルの精度を競う コンペティションです。予測結果を CSV ファイルに書き出して提出すると、順位が表示さ れます。(Kaggle に無料アカウント登録を行うと参加できます)

#### ▽ Write CSV オペレータを使用し、予測結果を CSV ファイルに書き出します



<sup>&</sup>lt;パラメータの設定>

パラメータ ×	環境 ×
🗳 Write CSV	
csv file	/prediction_SalePrice.csv
column separator	,
✓ write attribute nar	nes 🗊

フォルダのアイコンから保存場所を指定し、column separator はカンマ「,」に変更します。 プロセスを実行後、パラメータで指定した保存場所から CSV ファイルを開きます。 id と prediction(SalePreice)以外の列を削除し、prediction(SalePreice)の列名は SalePrice に 変更します。

	A	В	С	D
1	Id	SalePrice		
2	1461	119679		
3	1462	149361		
4	1463	177780		
-				

これで、コンペティションに参加する準備が完了しました。



▽ 予測結果 CSV ファイルを提出します

33			
Create	Overview Data Code Discussion	Leaderboard Rules Team	My Submissions Submit Predictions
Home	Step 1 Upload submission file		
			<b>↑</b>
Datasets			Upload Files
<> Code			
Discussions		File Format Your submission should be in CSV	Number of Predictions We expect the solution file to have 1459 prediction rows.
🕅 Courses		tormat. You can upload this in a zip/gz/rar/7z archive, if you prefer.	submission file on the data page.
✓ More			
ecently Viewed	Step 2 Describe submission	r ⊂ T B I	☞ " <> 📄 🖽 🛄 🙂 🖬
House Prices - Advanc		Briefly describe your submissio	200
<ul> <li>House Prices - Advanc</li> <li>(in Japanese) House P</li> </ul>		Briefly describe your submissio	n
<ul> <li>House Prices - Advanc</li> <li>(in Japanese) House P</li> <li>HPAR with Regression</li> </ul>		Briefly describe your submissio	n

Upload submission file から CSV ファイルをアップロードし、Make Submission をクリック すると提出完了です。順位表やスコアは、Leaderboard から確認できます。または、右上の アイコン(Your Profile>Competitions)からも順位が確認できます。

結果は、スコア 0.16167 で 4490 チーム中 2863 位の順位でした。



 House Prices - Advanced Regression Techniques
 2863/4490

 Predict sales prices and practice feature engineering, RFs, and gradient boosting
 Getting Started · Ongoing

### 終わりに

欠損値や外れ値など確認しながらデータ分析を行なっていきました。機械学習プロジェクトにおいて約8割がデータ前処理(データ準備)に費やされていると言われており、各変数の特徴や変数同士の関係性などを把握する上でも重要なプロセスになります。

今回初めて Kaggle のコンペティションに参加しました。データ分析の勉強にもなり、一か らプロセスを考えていく楽しさを感じながら進めていくことができました。皆さんも参加 されてみてはいかがでしょうか。

次回は、今回作成した予測モデルの精度を上げていきたいと思います。