

HOUSE PRICE PREDICTION✓

Problem Statement

- { i) How Price of the house is dependent on various features.
 ii) The top features that influence the price of the house.

A US-based housing company named Surprise Housing has decided to enter the Australian market. The company uses data analytics to purchase houses at a price below their actual values and flip them at a higher price. For the same purpose, the company has collected a data set from the sale of houses in Australia.

The company is looking at prospective properties to buy to enter the market.

2:17 PM Sat 19 Dec 64%

Assignment upGrad

Steps to proceed with the assignment.

1. Perform EDA to understand the data.
2. Check missing values.
3. For continuous variables, try to impute the missing values with mean or median. Perform EDA to find out which one fits best.
4. For categorical columns.
 - a. Check all categories carefully. → For those categorical columns, where you have meaningful missings, you should not use mode approach over here.

C_y
 A 90% } Skewed categorical columns
 B 5%
 C 3%
 D 2%

Basement quality
 1 1
 2 2
 3 3 "NO Basement"
 → NA: "NO Basement"

→ Categorical columns that are highly skewed

Activate Windows
Go to PC settings to activate Windows.

10

Try to drop skewed categorical models.

How to solve..Continue. upGrad

What to do?

5. Check if the target variable is normally distributed or not?
6. Create dummies for categorical data.
 - You can create groups of the the categories to reduce the number of categories and then create dummies.
 - This is an optional method.
7. Handling year columns.
 - There are 4 columns that contain year. What to do with them?
 - How to convert them?

Activate Windows
Go to PC settings to activate Windows.

11

Not removing outliers will influence the best fit line.

How to handle outliers?

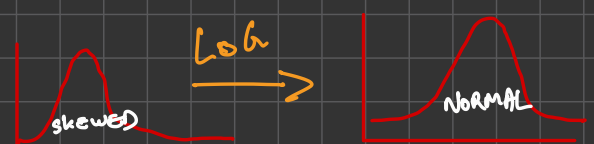
1- Do Log transformation for Sale Price

Linear Regression cannot extrapolate

- if $100 < y < 1000$
then it will always predict between $100 < \hat{y} < 1000$

- should not have outliers
→ capping
→ Drop outliers] not acceptable for assignment

→ Model won't work for range of data that was considered as outliers



Categorical Variable

- Group categories with lesser y of data into others

→ combine these categories to 'Rare' category

YEAR COLUMNS

- 1- calculate age from Yearbuilt column
- 2- use as continuous variable
- 3- Drop other Year features

MODELING

- 1- Start with RFE(50)
- 2- P-value/VIF
- 3- Ridge/Lasso } Tune hyperparameter alpha
- 4- Return top features

How to solve..Continue. upGrad

→ RFE (50) → P-value/VIF → Ridge / Lasso } hyperparameter tuning (λ / α)

What to do?

Model Building:

1. You can directly start trying out Lasso and Ridge with different values of alpha.
2. For lasso choose the best alpha.
3. For Ridge choose the best alpha.

After choosing the best alpha from both the models, check the performance of both the models.

Lastly, you need to find out best features that describes the price of the house, for this check the top features for both final models created using Ridge and Lasso and then choose the features accordingly.

Return the top features from both models

Activate Windows
Go to PC settings to activate Windows.

12

Lasso → Top 5 features

↓
Remove from TRAIN

↓
RETRAIN MODEL

↓
MEASURE PERFORMANCE

SUBJECTIVE QUESTIONS

I use code to answer in Jupyter

JUPYTER NOTEBOOK
PDF (SUBJECTIVE ANSWERS)

21P

3:05 PM Sat 19 Dec

×

Quiz Time

upGrad

Poll Questions

Question-1: Suppose we have a regularized linear regression model: $\arg\min ||Y - \beta x||^2 + \lambda ||\beta||$. What is the effect of increasing λ on bias and variance?

- ☐ a. Increases bias, increases variance
- ☒ b. Increases bias, decreases variance
- ☐ c. Decreases bias, increases variance
- ☐ d. Decreases bias, decreases variance
- ☐ e. Not enough information to tell

Question-2: After applying a regularization penalty in linear regression you find that some of the coefficients are zeroed out. Which of the following penalties might have been used?

- ☐ a. L0 norm
- ☒ b. L1 norm
- ☐ c. L2 norm
- ☐ d. either (A) or (B)
- ☐ e. any of the above

3:09 PM Sat 19 Dec

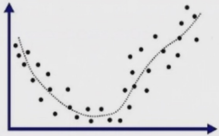
×

Quiz Time

upGrad

Poll Questions

Question-3: Suppose you used a degree-3 polynomial to fit the data and it look like as shown in the image, what will happen if we use a degree-4 polynomial?



- ☒ a. There are high chances that degree 4 polynomial will over fit the data
- ☐ b. There are high chances that degree 4 polynomial will under fit the data
- ☐ c. Can't say
- ☐ d. None of these

Question-4: Which of the following is true for the given fit (refer to the image)?

- ☐ a. Bias will be high, variance will be high
- ☐ b. Bias will be low, variance will be high
- ☐ c. Bias will be high, variance will be low
- ☒ d. Bias will be low, variance will be low

3:13 PM Sat 19 Dec

×

Quiz Time

upGrad

Poll Questions

Question-5: Suppose, you got a situation where you find that your linear regression model is under fitting the data. In such situation which of the following options would you consider?

- I will add more variables
- I will start introducing polynomial degree variables
- I will remove some variables

- ☒ a. 1 and 2
- ☐ b. 2 and 3
- ☐ c. 1 and 3
- ☐ d. 1, 2 and 3

Question-6: If the fitted model is underfitting then which of following regularization algorithm would you prefer?

- ☐ a. L1
- ☐ b. L2
- ☐ c. Any
- ☒ d. None of these

Handwritten notes: USED FOR AVOIDING OVERFITTING
Don't need to - SC

3:24 PM Sat 19 Dec

×

Quiz Time

upGrad

Poll Questions

Question-7: Let's say, a "Linear regression" model perfectly fits the training data (train error is zero). Now, Which of the following statement is true?

Diagram: Train data is perfectly fitted, but Test data shows some error.

- ☐ a. You will always have test error zero
- ☐ b. You can not have test error zero
- ☒ c. Can't say

Question-8: Suppose we fit "Lasso Regression" to a data set, which has 100 features (X_1, X_2, \dots, X_{100}). Now, we rescale one of these feature by multiplying with 10 (say that feature is X_1), and then refit Lasso regression with the same regularization parameter. Now, which of the following option will be correct?

Equation:
$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda p_1 + \lambda p_2 + \lambda p_3 + \dots + \lambda p_{100}$$
$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{100} x_{100}$$

Scale of β_1 depends on scale of x_1

- ☐ a. It is more likely for X_1 to be excluded from the model
- ☒ b. It is more likely for X_1 to be included in the model
- ☐ c. Can't say
- ☐ d. None of these

Handwritten notes: X_1 increases, β_1 decreases, whole model will not realize more



Quiz Time

upGrad

Poll Questions

Question-1: Suppose we have a regularized linear regression model: $\text{argmin} ||Y - \beta x||^2 + \lambda ||\beta||$. What is the effect of increasing λ on bias and variance?

- ☐ A • Increases bias, increases variance
- ☒ B • Increases bias, decreases variance
- ☐ C • Decreases bias, increases variance
- ☐ D • Decreases bias, decreases variance
- ☐ E • Not enough information to tell

Question-2: After applying a regularization penalty in linear regression you find that some of the coefficients are zeroed out. Which of the following penalties might have been used?

- ☐ A • L0 norm
- ☒ B • L1 norm
- ☐ C • L2 norm
- ☐ D • either (A) or (B)
- ☐ E • any of the above



Activate Windows
Go to PC settings to activate Windows.