

REGRESSION Project Write-up

Predicting house price in different cities and states of the USA

Marjan Rezvani

Abstract

House Price prediction is important to drive Real Estate efficiency, also customers, companies such as Zillow, and Banks can use this analysis to get better informed and decide marketing strategy.

The aim of this project is to predict the house prices using Machine learning algorithm Linear Regression considering factors such as lot square feet, number of bedrooms, number of bathrooms, address, city, state, zipcode and region.

I focused on real estate sales price in different cities of The United States, as listed on the Century21 website, and scraped the website using BeautifulSoup library. Then leveraged numerical and categorical feature engineering along with a linear regression to achieve promising results for this problem.

Design

First, I scraped the site of Century21 and extract useful features and then I did regularization and feature engineering on data and analyzed house price predications on different parameters and evaluated my model to get the best result that can best predict home prices given reasonable test/train splits in the data. I used linear regression with K-fold cross validation to predict the Target using the features listed below.

Data

The dataset I made after web scraping, contains 1753 observations and 10 attributes (9 predictors and 1 response). There were several variables that might help determine the sale price of houses such as lot square feet, number of bedrooms, number of bathrooms, address, city, state, and zipcode of the area.

Algorithms

Web Scraping:

Data is collected for various states and their cities from different regions by scraping the Century21 website for houses price.

Feature Engineering

1. Converting categorical features to binary dummy variables
2. Selecting subsets of the total unique values for categorical features that were converted to dummies
3. Converting data into a format required for a number of information processing needs by encoding
4. Doing particular calculations on some features to use them for modeling during EDA
5. Making new variables like region

Analysis:

Seaborn and matplotlib were used to do visualization of the data. The visualization gave some very good information about the trends we see with respect to the type of property.

Models

Machine Learning Linear Regression Model and Polynomial regression were used to find the model with strongest cross-validation performance.

Also tried Lasso and Ridge models, to get a model that minimizes errors

Model Evaluation and Selection

The entire training dataset was split into train, val, and test, and all scores reported below were calculated with 5-fold cross validation on the training portion only.

Linear Regression: 0.8437278269210082

Polynomial Regression: 0.9023538924173838

Lasso Regression: 0.7847769127129868

Ridge Regression: 0.712747428001326

Tools

- BeautifulSoup and Requests for web scraping and parsing the data
- NumPy and Pandas for data manipulation
- Scikit-learn for modeling and evaluation
- Matplotlib and Seaborn for plotting

Communication







