

DS203: Programming for Data Sciences

Assignment: Regression

Goal: This assignment aims to help you learn the application of Linear Regression for real estate evaluation. It involves understanding what data means, how to handle data, training the model, prediction, and testing your model. We will try to do the complete flow in this assignment.

Dataset: The market historical data set of real estate valuation is collected from Sindian District, New Taipei City, Taiwan. The dataset can be downloaded from [here](#).

Features Information:

The features of the dataset are as follows

X_1 = the transaction date (for example, 2013.250=2013 March, 2013.500=2013 June, etc.)

X_2 = the house age (unit: year)

X_3 = the distance to the nearest MRT station (unit: meter)

X_4 = the number of convenience stores in the living circle on foot (integer)

X_5 = the geographic coordinate, latitude. (unit: degree)

X_6 = the geographic coordinate, longitude. (unit: degree)

The output is as follow

Y = house price of unit area (3000 New Taiwan Dollar/meter squared)

More details about the dataset can be found in this [web-page](#).

What to do?

1. Download the dataset from above shared link.
2. Load the dataset into your python program and do pre-processing (remove the first row and first column as they are not useful).
3. Split loaded dataset into train and test dataset by keeping 80% samples in train dataset and remaining 20% samples in test dataset.
4. Now train linear regression model on train dataset (note that the last column (house price) is the output).
5. Report coefficients (weights corresponding to features) and intercept of trained model.
6. Predict price for every house (sample) in test dataset.
7. Compute mean squared error and R^2 value using predicted price and true price.
8. Repeat Step 4-6 for following train and test split: 60:40, 70:30, and 90:10. Report mean squared error and R^2 value for each split.
9. Use Ridge regression and Lasso models with following λ values (regularization parameter): 0.001, 0.005, 0.01, 0.05, 0.1, and 0.5. Report mean squared error and R^2 value for each alpha value with all train and test split ratios given in Step 7.

Practice on another dataset: (Assignment) You can use linear regression, Ridge regression, and LASSO models for predicting the heating load and cooling load requirements of buildings (that is, energy efficiency) as a function of building parameters. The energy efficiency dataset can download from [here](#) and more details about the dataset can be found in this [web-page](#).