

## Project Report

### Homework 2: Exploratory Data Analysis in IPython

This report discusses the different models used to solve the Zillows Challenge for predicting the log error of different properties on sale.

The data set consisted of the Properties file describing all the various properties about the parcel on sale which Zillow uses to predict the sales price. This data consisted of categorical, discrete and continuous values.

#### Following are the models used:

##### *1. Linear Regression Model*

Linear Regression follows a linear approach in modelling the relationship between scalar dependent variable 'Y' and one or more independent variables denoted by 'X'.

$$Y = mX + e$$

Where 'm' is the linear coefficient array which is multiplied to the X data set. And 'e' is the noise or induced error value.

Linear relationships are the simplest relationship that can be imagined between any 2 variables. Their simplicity makes it widely used across multiple applications in Statistics, Biology etc. Even if the variables are not linear, we can apply some transformation and reduce the gap between them.

##### *2. Random Forest Regression Model*

Random Forest model is part of ensemble learning methods which works by constructing multiple decision trees and outputting the class which is the mode or mean prediction of the decision tree. Random Forest Regression is a powerful tool capable of delivering performance that is most accurate to date. There are three additional built-in features of Random forest: performance assessment, measure of relative importance of descriptors and measure of compound similarity which is weighted by the relative importance of descriptors.

##### *3. K Nearest Neighbour Model*

KNN is a simple model which predicts the target based on similarity measure eg distance functions. A simple implementation of KNN regression calculates average of the numerical target based on the K nearest neighbours. For categorical variables KNN uses Euclidean distance and Manhattan distance to find the nearest neighbours.

$$\text{Euclidean} = \sqrt{\text{Sum of } (x_i - y_i)^2}$$

$$\text{Manhattan} = \text{Sum } |x_i - y_i|$$

Where  $x_i$  are all the training features and  $y_i$  corresponds to the target features.

## Analysis

In the Zillows problem surprisingly Linear Regression, Random Forest and KNN all gave similar mean square error, however compared to the coefficients in Linear the Feature importance values in Regression were a better estimation on the relation of these features with the log error.

Some of the attributes were highlighted in both the models as the most important ones like totalarea, taxamount, bedroomcnt. Thus concluding that these variables play a major role in determining the log error values.

Also one more interesting relationship that Random Forest model predicted was between latitude, longitude and log error. Thus suggesting that a parcel's location also plays a vital role in the selling price.

Also the Random Forest Regression model performed better when the features in the training data were reduced after the results of the previous regression values where most of the features were part of the training model.

However in spite of being powerful both Random and KNN model was not able to find the sample log error values within an accepted time frame and with the same data set columns Linear regression model was very quick to predict the log values. Thus the log errors predicted with Linear Regression model got a score of **0.0650533** and ranked **2344** on the Zillows Challenge

Featured Prediction Competition

**Zillow Prize: Zillow's Home Value Prediction (Zestimate)**  
Can you improve the algorithm that changed the world of real estate?

**\$1,200,000**  
Prize Money

Zillow · 3,368 teams · 4 months to go (5 days to go until merger deadline)

OverviewDataKernelsDiscussionLeaderboardRulesTeamMy SubmissionsSubmit Predictions

Your most recent submission

Name	Submitted	Wait time	Execution time	Score
sample_submission.csv	16 minutes ago	47 seconds	40 seconds	0.0650533

Complete

Jump to your position on the leaderboard ▾

Make a submission for **Zenab Bhinderwala**

2344newZenab Bhinderwala

0.0650533119m

Your Best Entry ↑

Your submission scored 0.0650533, which is not an improvement of your best score. Keep trying!