
CSE 519 -- Data Science (Fall 2017)
Homework 2: Exploratory Data Analysis in IPython

Name : Hemant Kumar Pandey
ID : 110828730

Report

27th September 2017

OVERVIEW

This homework is based on Kaggle challenge which involves a dataset by Zillow and revolves around predicting the price that a particular real estate property (usually a home) will sell for. The goal is to explore the data and uncover interesting observations about the Zillow data, model the data and submitting the results on Kaggle.

DATA PREPROCESSING

1. Dropping all the columns which have string values and does not contribute much to the correlation.
2. Conversion of the data types of all columns from object to numeric so that they can be manipulated faster and gives a better correlation.
3. Filling all the NAN values with the mean of respective column.
4. Standard scalar in the improvised model.

RANDOM FOREST REGRESSOR MODEL

How it works ?

It can be used for both classification and regression kind of problems. As the name suggest, this algorithm creates the forest with a number of trees. In general, the more rees in the forest the more robust the forest will be. In the same way in the random forest classifier, the higher the number of trees in the forest, we get more accurate results.

Random Forest is based on the decision tree concept which is a rule based system. If a training dataset is given with features and targets, the decision tree algorithm will come up with a set of rules. These rules can be used to make prediction on test dataset.

Suppose we want to see if an apartment will be sold for more than 200,000\$ or not ? To model the decision tree you will use the training dataset like the list of apartments in the nearby regions that were sold with the comparative price in the past.

So once you pass the dataset with the target as the property will be sold for more than 200,000\$ or not to the decision tree classifier. The decision tree will start building the rules with the apartments that were sold for good price as **nodes** and the targets which are not as the leaf nodes. By considering the path from the root node to the leaf node, we can get the rules. The simple rule could be if some apartment has more than 2 bedrooms and a pool, it will be sold for more than 200,000\$. We just need to check the rules which are created by the decision tree to predict the result.

In random forest algorithm, Instead of using information gain or gini index for calculating the root node, the process of finding the root node and splitting the feature nodes will happen randomly.

How well it works ?

- Handles the missing values.
- When there are more trees in forest, it would not overfit the model.
- Can be used for both classification and regression tasks.

INTERESTING THINGS I LEARNED

- First experience with modelling and Jupyter.
- How selecting different variables can change the result which can vary much.
- The power of pandas and sklearn. It was very easy to use.

KAGGLE SCORE FROM BOTH MODELS

Score from both models

Overview
Data
Kernels
Discussion
Leaderboard
Rules
Team
My Submissions
Submit Predictions

You can select up to 2 submissions to be used to calculate your final leaderboard score. If 2 submissions are not selected, they will be chosen based on your best submission scores on the public leaderboard.

Your final score may not be based on the same exact subset of data as the public leaderboard, but rather a different private data subset of your full submission — your public score is only a rough indication of what your final score is.

You should thus choose submissions that will most likely be best overall, and not necessarily on the public subset.

0 submissions for HemantPandey

Sort by Most recent

All

Successful

Selected

Submission and Description	Public Score	Use for Final Score
RandomForest_predicted_results.csv 4 minutes ago by HemantPandey add submission details	0.0650277	<input type="checkbox"/>
linearRegression_predicted_results.csv an hour ago by HemantPandey add submission details	0.0672694	<input type="checkbox"/>

Rank On Kaggle

Overview
Data
Kernels
Discussion
Leaderboard
Rules
Team
My Submissions
Submit Predictions

2277	▼ 43	Sid 2		0.0650247	18	4d
2278	▼ 43	everything.bigdata		0.0650276	30	3h
2279	new	HemantPandey		0.0650276	2	12m
<div> Your Best Entry ↗ </div> <div> You advanced 965 places on the leaderboard! Your submission scored 0.0650277, which is an improvement of your previous score of 0.0672694. Great job! </div> <div> Tweet this! </div>						
2280	new	GovindShukla		0.0650279	2	13m
2281	▼ 45	Yuxin Tian		0.0650280	3	2mo
2282	▼ 45	Alireza21		0.0650295	39	16d
2283	▼ 45	Aurel Iuga		0.0650298	12	2mo
2284	▼ 45	lianqiao		0.0650298	4	4mo
2285	▼ 45	Adit Deshpande		0.0650303	4	3mo
2286	▼ 45	RiddhiRexAntonyrex		0.0650305	5	14h
2287	▼ 45	aminos		0.0650320	3	1mo