

3. Exploring data tables with Pandas

1. Use Pandas to read the house prices data. How many columns and rows are there in this dataset?
2. The first step I usually do is to use commands like `pandas.head()` to print a few rows of data. Look around what kind of features are available and read `data description.txt` for more info. Try to understand as much as you can. Pick three features you think will be good predictors of house prices and explain what they are.
3. How many unique conditions are there in `SaleCondition`? Use Pandas to find out how many samples are labeled with each condition. What do you learn from doing this?
4. Select one variable you picked in b., do you want to know something more about that variable? Use Pandas to answer your own question and describe what you did shortly here.

In []:

4. Learning to explore data with Seaborn

1. Let us first look at the variable we want to predict `SalePrice`. Use Seaborn to plot histogram of sale prices. What do you notice in the histogram?
2. Plot the histogram of the `LotArea` variable. What do you notice in the histogram?
3. Use Seaborn to plot `LotArea` in the x-axis and `SalePrice` on the y-axis. Try plotting `log(LotArea)` versus `log(SalePrice)` and see if the plot looks better.

In []:

5. Dealing with missing values

1. Suppose we want to start the first step of house price modeling by exploring the relationship between four variables: `MSSubClass`, `LotArea`, `LotFrontage` and `SalePrice`. I have done some exploring and found out that `LotFrontage` has a lot of missing values, so you need to fix it.
2. `LotFrontage` is the width of the front side of the property. Use Pandas to find out how many of the houses in our database is missing `LotFrontage` value.
3. Use Pandas to replace `NaN` values with another number. Since we are just exploring and not modeling yet, you can simply replace `NaN` with zeros for now.

In []:

6. Correlations between multiple variables

One incredible feature of Seaborn is the ability to create correlation grid with `pairplot` function. We want to create one single plot that show us how all variables are correlated.

1. First, you need to create a data table with four columns: `MSSubClass`, `LotArea` (with `log` function applied), `LotFrontage` (missing values replaced) and `SalePrice` (with `log` function applied).
2. Then, use `pairplot` to create a grid of correlation plots. What do you observe from this plot?

In []:

7. Data Preparation

Let's prepare train.csv for model training

1. Pick columns that are numeric data and plot distributions of those data (with Seaborn). If you find a column with skewed distribution you will write a script to transform that column with a log function. Then standardize them.
2. For categorical variables, we will simply transform categorical data into numeric data by using function `pandas.get_dummies()`.
3. Split data into x and y. The variable x contains all the house features except the SalePrice. y contains only the SalePrice.

In []:

8. Let us first fit a very simple linear regression model, just to see what we get.

1. Use import LinearRegression from sklearn.linear model and use function `fit()` to fit the model.
2. Use function `predict()` to get house price predictions from the model (let's call the predicted house prices yhat).
3. Plot y against yhat to see how good your predictions are.

In []:

9. Assessing Your Model

According to Kaggle's official rule on this problem, they use root mean square errors (rmse) to judge the accuracy of our model. This error computes the difference between the log of actual house prices and the log of predicted house price. Find the mean and squareroot them.

We want to see how we compare to other machine learning contestants on Kaggle so let us compute our rmse. Luckily, sklearn has done most of the work for you by providing mean square error function. You can use it by importing the function from sklearn.metrics. Then, you can compute mean square error and take a squareroot to get rmse.

What's the rmse of your current model? Check out Kaggle Leaderboard for this problem to see how your number measures up with the other contestants.

In []:

10. Cross Validation

As we discussed earlier, don't brag about your model's accuracy until you have performed cross validation. Let us check cross-validated performance to avoid embarrassment.

Luckily, scikit learn has done most of the work for us once again. You can use the function `cross_val_predict()` to train the model with cross validation method and output the predictions.

What's the rmse of your cross-validated model? Discuss what you observe in your results here. You may try plotting this new yhat with y to get better insights about this question.

In []:

11 (Optional) Fit Better Models

There are other models you can fit that will perform better than linear regression. For example, you can fit linear regression with L2 regularization. This class of models has a street name of 'Ridge Regression' and sklearn simply called them Ridge. As we learned last time, this model will fight overfitting problem.

Furthermore, you can try linear regression with L1 regularization (street name Lasso Regression or Lasso in sklearn). Try these models and see how you compare with other Kagglers now. You can write about your findings below.

In []: