

# Kaggle Quick Start

An incomplete Starter Kit for House Price Prediction  
with Regression

Woodmark

## 1. Quick Start in Python or R

	Python	R
Import necessary libraries	<pre>import numpy as np import pandas as pd import matplotlib.pyplot as plt from scipy.stats import skew from sklearn import linear_model from sklearn.preprocessing import StandardScaler from sklearn.model_selection import cross_val_score</pre>	<pre>install.packages('moments') install.packages('glmnet')  library(moments) library(glmnet)</pre>
Load input data	<pre>train = pd.read_csv("input/train.csv") test = pd.read_csv("input/test.csv")</pre>	<pre>getwd() setwd() train &lt;- read.csv("train.csv") test &lt;- read.csv("test.csv")</pre>
Check data dimensions	<pre>print("Number of features in training set: {}".format(train.shape[1])) print("Number of training data entries: {}".format(train.shape[0])) print("Number of test data entries: {}".format(test.shape[0]))</pre>	<pre>dim(train) dim(test) str(train) str(test) summary(train) summary(test)</pre>
Where is the ID / the training prices?	<pre>print("First column in both sets is: {}".format(train.columns[0])) print("Last column in training set is: {}".format(train.columns[-1]))</pre>	-
Make index column the actual data frame index	<pre>train.set_index('Id', inplace=True) test.set_index('Id', inplace=True)</pre>	-
Split off the price column from the training data:	<pre>train_price = train["SalePrice"] train.drop("SalePrice", axis=1, inplace=True)</pre>	<pre>train_price&lt;- train\$SalePrice train\$SalePrice &lt;- NULL</pre>
Check value ranges of features variables:	<pre>%matplotlib inline plt.rcParams['figure.figsize'] = (16, 6) train.boxplot(showfliers=False, rot=90) plt.show()</pre>	<pre>nums &lt;- sapply(train, is.numeric) boxplot(train[,nums], use.cols=TRUE, las = 2)</pre>
Deal with numerical features:	<pre># extract locations of numerical features num_feat = (train.dtypes != "object").as_matrix()  train_num = train.iloc[:, num_feat] test_num = test.iloc[:, num_feat]</pre>	<pre>num_feat &lt;- sapply(train, is.numeric) num_feat2&lt;- sapply(test, is.numeric)  train_num&lt;-train[,num_feat] test_num&lt;-test[,num_feat2]</pre>
Check positivity and fill missing values with	<pre>print("All numerical values in training set positive? {}".format(not (train_num &lt; 0).any().any()))</pre>	<pre>which(train_num&lt;0) which(test_num&lt;0) for(i in 1:ncol(train_num)){   train_num[is.na(train_num[,i]), i] &lt;-</pre>

column means:	<pre># fill missing values in training set with column means train_num = train_num.fillna(train_num.mean()) train.iloc[:, num_feat] = train_num  print("All numerical values in test set positive? {}".format(not (test_num &lt; 0).any().any())) # fill missing values in test set with TRAINING column means test_num = test_num.fillna(train_num.mean()) test.iloc[:, num_feat] = test_num</pre>	<pre>mean(train_num[i], na.rm = TRUE))  for(i in 1:ncol(test_num)){   test_num[is.na(test_num[,i]), i] &lt;- mean(train_num[,i], na.rm = TRUE)}</pre>
Check skewness	<pre>ax = train_price.plot.density() ax.set_xlim([0,train_price.max()]) ax.set_xlabel("price") ax.set_ylabel("density") plt.show()</pre>	<pre>skewness(train_num)</pre>
	<pre>print("Numerical feature columns: {}".format(train.columns[num_feat])) print("Skeweness of numerical training features: {}".format(skew(train_num)))</pre>	-
log(1+p) transform of skewed features	<pre>skewed = (np.absolute(skew(train_num)) &gt; 1) train_num.iloc[:, skewed] = np.log1p(train_num.iloc[:, skewed]) test_num.iloc[:, skewed] = np.log1p(test_num.iloc[:, skewed]) train_price = np.log1p(train_price)</pre> <pre>print("Skeweness of numerical training features after transformation: {}".format(skew(train_num)))</pre>	<pre>for (col in colnames(train_num)){ if(abs (skewness(train_num[,col]))&gt;1){ train_ num[,col]&lt;-log(1+train_num[,col]) }  for (col in colnames(test_num)){ if(abs( skewness(test_num[,col]))&gt;1){ test_nu m[,col]&lt;-log(1+test_num[,col]) }</pre>
Normalize numerical features	<pre>scaler = StandardScaler().fit(train_num) train_num = scaler.trans- form(train_num) test_num = scaler.transform(test_num</pre>	<pre>for (col in colnames(train_num)){ train_ num[,col]&lt;-scale(train_num[,col]) }  for (col in colnames(test_num)){ test_ num[,col]&lt;-scale(test_num[,col]) }</pre>
Apply the transformed values to the original sets	<pre>train.iloc[:, num_feat] = train_num test.iloc[:, num_feat] = test_num</pre>	<pre>train[,num_feat]&lt;-train_num test[,num_feat2]&lt;-test_num</pre>
Transform categorical features	<pre>train_test = pd.concat([train, test]) train_test = pd.get_dummies(train_test) train_test.head()</pre>	...
Split sets again	<pre>train = train_test.iloc[:train.shape[0],:] test = train_test.iloc[train.shape[0]:,:]</pre>	...

Linear Regression	<pre>lin_reg = linear_model.LinearRegression() scores = cross_val_score(lin_reg, train,                           train_price,                           cv=5, scoring='neg_mean_squared_error') print("Mean of 5 CV sqrt MSE: {}".format(np.sqrt(-scores.mean())))</pre>	...
Improve the Model	...	...

## 2. Ressources

- Kaggle tutorial links: <https://www.kaggle.com/wiki/Tutorials>
- Scikit-Learn Python Machine Learning Library: <http://scikit-learn.org/>
- Tools:
  - R Studio <https://www.rstudio.com/>
  - Anaconda <https://www.continuum.io/downloads>
- Statistics:
  - Statistik: Weg zur Datenanalyse (Ludwig Fahrmeir)
  - Theoretical Statistics: Topics for a Core Course (Robert W. Keener)
- Probability Theory:
  - Wahrscheinlichkeitstheorie (Achim Klenke)
  - Understanding Probability (Henk Tijms)
- Machine Learning:
  - An Introduction to Statistical Learning <http://www-bcf.usc.edu/~gareth/ISL/>
  - The Elements of Statistical Learning <http://statweb.stanford.edu/~tibs/Elem-StatLearn/>
- Online Courses:
  - Coursera Data Science Specialization <https://www.coursera.org/specializations/jhu-data-science>
  - Coursera Machine Learning <https://www.coursera.org/learn/machine-learning>
- Cheat Sheets
 

R:

  - <https://www.rstudio.com/wp-content/uploads/2016/10/r-cheat-sheet-3.pdf>
  - <https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheat-sheet.pdf>

- [http://nicolascampione.weebly.com/uploads/1/9/4/1/19411255/r\\_cheat\\_sheet.pdf](http://nicolascampione.weebly.com/uploads/1/9/4/1/19411255/r_cheat_sheet.pdf)

Python:

- [https://perso.limsi.fr/pointal/\\_media/python:cours:mementopython3-english.pdf](https://perso.limsi.fr/pointal/_media/python:cours:mementopython3-english.pdf)
- <https://ehmatthes.github.io/pcc/cheatsheets/README.html>
- [https://s3.amazonaws.com/assets.datacamp.com/blog\\_assets/PythonForDataScience.pdf](https://s3.amazonaws.com/assets.datacamp.com/blog_assets/PythonForDataScience.pdf)

Regression:

- <http://www.alisonpearce.net/wp-content/uploads/2013/05/Fun-Reg-Cheat-Sheet.pdf>
- <http://www.stat.ucla.edu/~rosario/classes/091/112-1b/regression>