

Regression Week 1: Simple Linear Regression Assignment

Predicting House Prices (One feature)

In this notebook we will use data on house sales in King County, where Seattle is located, to predict house prices using simple (one feature) linear regression. You will:

- Use SArray and SFrame functions to compute important summary statistics
- Write a function to compute the Simple Linear Regression weights using the closed form solution
- Write a function to make predictions of the output given the input feature
- Turn the regression around to predict the input/feature given the output
- Compare two different models for predicting house prices

If you are doing the assignment with IPython Notebook

An IPython Notebook has been provided below to you for this assignment. This notebook contains the instructions, quiz questions and partially-completed code for you to use as well as some cells to test your code.

What you need to download

If you are using GraphLab Create:

- Download the King County House Sales data in SFrame format: `kc_house_data.gl.zip`
(https://eventing.coursera.org/api/redirectStrict/UbzBRAX6rjRm88u-rNTtjg3fywGHJ06mKp000In0PdZ5qrQ909aGtMjNIAKid-C2gSeZxd4AwZxj58dzoBj01Q.HEJtxnM3fEt10nuEEM7w8A.6L63RV-_LKzMKbYU8_xBRvYAVVjnFdEQCou5rZPsfvWlwmaz6ipuoFKTznqVR7rz8dwCjrffZxX45sOKWMAioR4yGz6OJPPXauJrgB9VV4GUhtpCIEYlleTQ_57iq3gsyjjGr8S7s3gLo039FN9E_3K7XtcODTK8ipzsW5vH12SHiZO7ny3XEX9F1efEpmpiqOv8Dcaoa74sTs2YktCOWkcd4yC9cUAX_2k_RC44h7LGQGljLow7T1gsOb5bXkOlsh1Tej64EYviM-Oj0gRz_pdENYgQz6kxKSv3kDQrK6-ildrQah-KajpEojl7GETefn1plowLVeAHllsYl9s7yuyFsBhUuP5JsPwNGF4nA6mJ7nzhAautl-kLq0XEwLirUdvpDZFs78Tb0S1dCGhKBOBiXPFa9AhrZmQgoRK3MXazbjAGHajTj5deIdclIOYll)
- Download the companion IPython Notebook: `week-1-simple-regression-quiz-blank.ipynb`
(https://eventing.coursera.org/api/redirectStrict/AeCRGpGCz3P2-H2-6i3g2hsCRp38npEpK53SOW-oWZiKzX8_9vDh5j1N0-KjqBl_bC1RWaHpb76Ali1S5k0BQg.NzbPV1bszGwhDlug6lzk1Q.GcXfG1FrZUfz-cffkT5vGQCqLqK8anbjenV8Zp3cMX7T6PnGbojNUJnFsj1QECp04RodYkDAUhl9x2ov6dm4zCcQdJeYW_YAyyj-Ign9Cr_JpEomzpKlqZ9RQoXXQPO-mPaYsrHbGZwjNFSRD0ajc-KWPlomJFGKMqG7zk_YEba-sbCDxYaNvDgKB8iY9C5PO60QAQ-gN_4pyA5otpPE5V69Ka1GbOeEugoHw3F8w7UoMj8PZ1hv4FoVt7R6VF_9zdc4-udQAdnSWVGzZ8GOBCKrLQtKU1tpAjS0FwjVW_x9jtU_tqRyv-uSZG9TLpb7INPQvXRunCKs48SOPnlHg3uDvRMXiQcFO3vGAhQxrwhYl0-NusbK5tqWDKHeF8c0EzBYA4-JzekdzQK6wDZdq3fpWwSsLRWeCzWvP_4BP2XP8RyCvzvOLPIP5T6pAx0PR3avD7dYWPpERHsw2F8kCxcg5fMjdr51fkHmQXvQF5LFtHcdXON5rwXICFzHFSly_3z)

- Save both of these files in the same directory (where you are calling IPython notebook from) and unzip the data file.

If you are not using GraphLab Create:

- Download the King County House Sales data csv file: `kc_house_data.csv` (https://eventing.coursera.org/api/redirectStrict/85fSP95hXo3_SDqHAeRcC2xkePX-LWSSNhL3b7zb3RunesPA-aolUWypNlKmch0Zj1VXuKi4GaGsU0v4-tZmqQ.8MLZu6Q39vdsOe8psePwTw.C_-kPAxEDcIA2rA_52z6CRsDXW2bGANw0Gmm0qaiUU-MbcjMrbFxF4Nv4CUjt7FwRirI6TI8h0FvZsxdQOJNlYZM_bxs4cNi1mUCVeQQJcpjZDFK-sgl6758VuVMIUsxU65OBPwq5ScUOmFjyO0sOCpEWy_-EpSVjI7DKEyzBY1GAX9cqej4pa3my9oiyxPsf-j6mNeRhF3FvOZHdzuzdUs9lZAZkcsNWblvbfh_Tf--NNv4ywscjRPTXPPn0qWfEs3ONZSjPpxnOvHKae4HuB8s1OLPH1JDs8p1SfRFGib3tpfl6W7FGCiOo_hnxThavNI4dj0zcOqQrYy-hgVIOvj55qt9BYp4l223nEck3Voj9X6EkgbTh8E70ui6KvcALMxFf86W62gapmaSIS778JG8uVE489vMd_Li-HBgVEz5xP8ft0Yw-NpfeRPNEO)
- Download the King County House Sales training data csv file: `kc_house_train_data.csv` (https://eventing.coursera.org/api/redirectStrict/KjEV4A2hdDfwfrKfVM4D3jy4ni0MfzF2090UnhJdRa5xUd9nWp-XOH3eDY1vNHYoZLiRjBefFI4BEH11ojj5Q.xztW974TR9jxihm3DR-eaQ.Illcoi0eXLTMH0zqPdIj6uejc5O6FLRWCDtyGkGbvjNlnvj0KA87AmTbaVqcDPZ9XmZDsYdl814rli0vzfPlcyzhN5qoQc3_HAShGHE1Ek6WtZt-4CZNPiw2lO4LyMbg-NY9FZpMUKQiz5-3pHeRiywToHkyy4ZasyYTQsb9TX0yEORvvpd9l9mvmmbambVk_-nrUPnhdbGqCm3r_w6F7KbwePBz2cpkb9SV-RnfYIV4_IBRUymVhqw5u5hEqYaGaO5bKQC4vcuqHgGxExKj54wxsdC4v02wbLuibDpfQJZZR9ZCQ2NBL51PvjoNPcbz9QceVH3kpDO6cdwMPP1JDeTGeLB6_GVygXZh8Yx-Zq_Imlkz1TbrXWO3JT-3HhgZ7EsdqetTjJW346Q7_anbD8WDIXeVKAcb6CDmiOd4LFpXgq6tNFEDrviRPb8mckiVqoILnagc75GidyelAMNln5A)
- Download the King County House Sales testing data csv file: `kc_house_test_data.csv` (https://eventing.coursera.org/api/redirectStrict/rAvkNMxf72p_FhvdOv7djw8lb-kDy_XvBCFW5PZdgU_bevxCoDxXs7fmXUI4Plvi7O738sjNrtj993D4TGcWcw.Sm55cYcTIEguUcu0mmCVUw.-0q0STAcYl9bGOwoaevT6x-oeK7SsE-zzrwMzksoixQ8JbEO9LI7GKyqMwAMnUQQ9lRgOuGXMMYLR55_RkPu0iWE_jjbZIP9qiSl9P8YnEF6EsxwRpMc6cNlxIcjcZC1bVg0lj1_6DaVRFGW507mrIWyf-uLjDlFFCjCiSxgQxC88f2bMV5-TuskYllsHrxHjyU6wSHS861rxnSbKQCKtLFG2bLy9khKFHWc4UXUDyvQReCGpJckIXgoJmw3uNXpSthTURV-k7lXWGcA8b77s9reZm9bm7FOq43FV5xqMSr3Av2nEj1dTC41RbuyN2kS5a8vZm8Q7-A0shYRbVm1XK0a9PmNuH8jRy1jVMj1Oig-6zUA_VwNkqHIUO2OscpoAVs0fi0y_nBASucmKit_fKCs1qeNua_wjb5R7P38AayMrhO-QhvVnJ84WYyIK3)

Useful resources

You may need to install the software tools or use the free Amazon EC2 machine. Instructions for both options are provided in the reading for Module 1.

If instead you are using other tools to do your homework

You are welcome, however, to write your own code and use any other libraries, like Pandas or R, to help you in the process. If you would like to take this path, follow the instructions below.

1. If you are using SFrame, import graphlab and load in the house data, otherwise you can also download the csv. (Note that we will be using the training and testing csv files provided). e.g in python with SFrames:

```
sales = graphlab.SFrame('kc_house_data.gl/')
```

2. Split data into 80% training and 20% test data. Using SFrame, use this command to set the same seed for

everyone. e.g. in python with SFrames:

```
train_data, test_data = sales.random_split(.8, seed=0)
```

For those students not using graphlab please download the training and testing data csv files.

From now on we will train the models using train_data. It will be important that we use the same split here to ensure the results are the same.

3. Write a generic function that accepts a column of data (e.g. an SArray) 'input_feature' and another column 'output' and returns the Simple Linear Regression parameters 'intercept' and 'slope'. Use the closed form solution from lecture to calculate the slope and intercept. e.g. in python:

```
def simple_linear_regression(input_feature, output):
    [your code here]
    return(intercept, slope)
```

4. Use your function to calculate the estimated slope and intercept on the training data to predict 'price' given 'sqft_living'. e.g. in python with SFrames using:

```
input_feature = train_data['sqft_living']
output = train_data['price']
```

save the value of the slope and intercept for later (you might want to call them e.g. squarefeet_slope, and squarefeet_intercept)

5. Write a function that accepts a column of data 'input_feature', the 'slope', and the 'intercept' you learned, and returns an a column of predictions 'predicted_output' for each entry in the input column. e.g. in python:

```
def get_regression_predictions(input_feature, intercept, slope)
    [your code here]
    return(predicted_output)
```

6. Quiz Question: Using your Slope and Intercept from (4), What is the predicted price for a house with 2650 sqft?

7. Write a function that accepts column of data: 'input_feature', and 'output' and the regression parameters 'slope' and 'intercept' and outputs the Residual Sum of Squares (RSS). e.g. in python:

```
def get_residual_sum_of_squares(input_feature, output, intercept, slope):
    [your code here]
    return(RSS)
```

Recall that the RSS is the sum of the squares of the prediction errors (difference between output and prediction).

8. Quiz Question: According to this function and the slope and intercept from (4) What is the RSS for the simple linear regression using squarefeet to predict prices on TRAINING data?

9. Note that although we estimated the regression slope and intercept in order to predict the output from the input, since this is a simple linear relationship with only two variables we can invert the linear function to estimate the input given the output!

Write a function that accept a column of data: 'output' and the regression parameters 'slope' and 'intercept' and outputs the column of data: 'estimated_input'. Do this by solving the linear function $\text{output} = \text{intercept} + \text{slope} * \text{input}$ for the 'input' variable (i.e. 'input' should be on one side of the equals sign by itself). e.g. in python:

```
def inverse_regression_predictions(output, intercept, slope):  
    [your code here]  
    return(estimated_input)
```

10. Quiz Question: According to this function and the regression slope and intercept from (3) what is the estimated square-feet for a house costing \$800,000?
11. Instead of using 'sqft_living' to estimate prices we could use 'bedrooms' (a count of the number of bedrooms in the house) to estimate prices. Using your function from (3) calculate the Simple Linear Regression slope and intercept for estimating price based on bedrooms. Save this slope and intercept for later (you might want to call them e.g. bedroom_slope, bedroom_intercept).
12. Now that we have 2 different models compute the RSS from BOTH models on TEST data.
13. Quiz Question: Which model (square feet or bedrooms) has lowest RSS on TEST data? Think about why this might be the case.