

Regression Week 2: Multiple Regression (Interpretation)

The goal of this first notebook is to explore multiple regression and feature engineering with existing graphlab functions.

In this notebook you will use data on house sales in King County to predict prices using multiple regression. You will:

- Use SFrames to do some feature engineering
- Use built-in graphlab functions to compute the regression weights (coefficients/parameters)
- Given the regression weights, predictors and outcome write a function to compute the Residual Sum of Squares
- Look at coefficients and interpret their meanings
- Evaluate multiple models via RSS

Fire up graphlab create

```
In [72]:  
import graphlab
```

Load in house sales data

Dataset is from house sales in King County, the region where the city of Seattle, WA is located.

```
In [73]:  
sales = graphlab.SFrame('kc_house_data.gl/')
```

Split data into training and testing

We use seed=0 so that everyone running this notebook gets the same results. In practice, you may set a random seed (or let GraphLab Create pick a random seed for you).

```
In [74]:  
train_data, test_data = sales.random_split(.8, seed=0)
```

Learning a multiple regression model

Recall we can use the following code to learn a multiple regression model predicting 'price' based on the following features: example_features = ['sqft_living', 'bedrooms', 'bathrooms'] on training data with the following code:

(Aside: We set validation_set = None to ensure that the results are always the same)

```
In [75]:  
example_features = ['sqft_living', 'bedrooms', 'bathrooms']  
example_model = graphlab.linear_regression.create(train_data, target =  
'price', features = example_features,  
                                                  validation_set = None)
```

PROGRESS: Linear regression:

PROGRESS: -----

PROGRESS: Number of examples : 17384

PROGRESS: Number of features : 3

```

PROGRESS: Number of unpacked features : 3
PROGRESS: Number of coefficients      : 4
PROGRESS: Starting Newton Method
PROGRESS: -----
PROGRESS: +-----+-----+-----+-----+-----+
-----+
PROGRESS: | Iteration | Passes   | Elapsed Time | Training-max_error |
Training-rmse |
PROGRESS: +-----+-----+-----+-----+-----+
-----+
PROGRESS: | 1          | 2        | 0.009000     | 4146407.600631     |
258679.804477 |
PROGRESS: +-----+-----+-----+-----+-----+
-----+
PROGRESS: SUCCESS: Optimal solution found.
PROGRESS:

```

Now that we have fitted the model we can extract the regression weights (coefficients) as an SFrame as follows:

```

In [76]:
example_weight_summary = example_model.get("coefficients")
print example_weight_summary

```

```

+-----+-----+-----+
|      name      | index |      value      |
+-----+-----+-----+
| (intercept)    | None  | 87910.0724924   |
| sqft_living     | None  | 315.403440552   |
| bedrooms       | None  | -65080.2155528  |
| bathrooms      | None  | 6944.02019265   |
+-----+-----+-----+
[4 rows x 3 columns]

```

Making Predictions

In the gradient descent notebook we use numpy to do our regression. In this book we will use existing graphlab create functions to analyze multiple regressions. Recall that once a model is built we can use the `.predict()` function to find the predicted values for data we pass. For example using the example model above:

```

In [77]:
example_predictions = example_model.predict(train_data)
print example_predictions[0] # should be 271789.505878

271789.505878

```

Compute RSS

Now that we can make predictions given the model, let's write a function to compute the RSS of the model. Complete the function below to calculate RSS given the model, data, and the outcome.

```

In [78]:
def get_residual_sum_of_squares(model, data, outcome):

```

```

# First get the predictions
predictions = model.predict(data)

# Then compute the residuals/errors
residual = outcome - predictions

# Then square and add them up
residual_squared = residual * residual
RSS = residual_squared.sum()

return(RSS)

```

Test your function by computing the RSS on TEST data for the example model:

```

In [79]:
rss_example_train = get_residual_sum_of_squares(example_model, test_data,
test_data['price'])
print rss_example_train # should be 2.7376153833e+14

2.7376153833e+14

```

Create some new features

Although we often think of multiple regression as including multiple different features (e.g. # of bedrooms, squarefeet, and # of bathrooms) but we can also consider transformations of existing features e.g. the log of the squarefeet or even "interaction" features such as the product of bedrooms and bathrooms.

You will use the logarithm function to create a new feature. so first you should import it from the math library.

```

In [80]:
from math import log

```

Next create the following 4 new features as column in both TEST and TRAIN data:

- bedrooms_squared = bedrooms*bedrooms
- bed_bath_rooms = bedrooms*bathrooms
- log_sqft_living = log(sqft_living)
- lat_plus_long = lat + long As an example here's the first one:

```

In [81]:
train_data['bedrooms_squared'] = train_data['bedrooms'].apply(lambda x:
x**2)
test_data['bedrooms_squared'] = test_data['bedrooms'].apply(lambda x:
x**2)
In [82]:
# create the remaining 3 features in both TEST and TRAIN data
train_data['bed_bath_rooms'] = train_data['bedrooms'] *
train_data['bathrooms']
test_data['bed_bath_rooms'] = test_data['bedrooms'] *
test_data['bathrooms']

```

```

train_data['log_sqft_living'] = train_data['sqft_living'].apply(lambda x:
log(x))
test_data['log_sqft_living'] = test_data['sqft_living'].apply(lambda x:

```

```
log(x))
```

```
train_data['lat_plus_long'] = train_data['lat'] + train_data['long']  
test_data['lat_plus_long'] = test_data['lat'] + test_data['long']
```

- Squaring bedrooms will increase the separation between not many bedrooms (e.g. 1) and lots of bedrooms (e.g. 4) since $1^2 = 1$ but $4^2 = 16$. Consequently this feature will mostly affect houses with many bedrooms.
- bedrooms times bathrooms gives what's called an "interaction" feature. It is large when *both* of them are large.
- Taking the log of squarefeet has the effect of bringing large values closer together and spreading out small values.
- Adding latitude to longitude is totally non-sensical but we will do it anyway (you'll see why)

Quiz Question: What is the mean (arithmetic average) value of your 4 new features on TEST data? (round to 2 digits)

```
In [83]:  
print(test_data['bedrooms_squared'].mean())  
print(test_data['bed_bath_rooms'].mean())  
print(test_data['log_sqft_living'].mean())  
print(test_data['lat_plus_long'].mean())
```

```
12.4466777016  
7.50390163159  
7.55027467965  
-74.6533349722
```

Learning Multiple Models

Now we will learn the weights for three (nested) models for predicting house prices. The first model will have the fewest features the second model will add one more feature and the third will add a few more:

- Model 1: squarefeet, # bedrooms, # bathrooms, latitude & longitude
- Model 2: add bedrooms*bathrooms
- Model 3: Add log squarefeet, bedrooms squared, and the (nonsensical) latitude + longitude

```
In [84]:  
model_1_features = ['sqft_living', 'bedrooms', 'bathrooms', 'lat', 'long']  
model_2_features = model_1_features + ['bed_bath_rooms']  
model_3_features = model_2_features + ['bedrooms_squared',  
    'log_sqft_living', 'lat_plus_long']
```

Now that you have the features, learn the weights for the three different models for predicting target = 'price' using `graphlab.linear_regression.create()` and look at the value of the weights/coefficients:

```
In [85]:  
# Learn the three models: (don't forget to set validation_set = None)  
model_1 = graphlab.linear_regression.create(train_data, target = 'price',  
    features = model_1_features,  
    validation_set = None)  
model_2 = graphlab.linear_regression.create(train_data, target = 'price',  
    features = model_2_features,  
    validation_set = None)
```

```
model_3 = graphlab.linear_regression.create(train_data, target = 'price',
features = model_3_features,
validation_set = None)
```

PROGRESS: Linear regression:

PROGRESS: -----

PROGRESS: Number of examples : 17384

PROGRESS: Number of features : 5

PROGRESS: Number of unpacked features : 5

PROGRESS: Number of coefficients : 6

PROGRESS: Starting Newton Method

PROGRESS: -----

PROGRESS: +-----+-----+-----+-----+-----+-----+-----+
-----+

PROGRESS:	Iteration		Passes		Elapsed Time		Training-max_error	
Training-rmse								

PROGRESS: +-----+-----+-----+-----+-----+-----+-----+
-----+

PROGRESS:	1		2		0.018001		4074878.213096	
236378.596455								

PROGRESS: +-----+-----+-----+-----+-----+-----+-----+
-----+

PROGRESS: SUCCESS: Optimal solution found.

PROGRESS:

PROGRESS: Linear regression:

PROGRESS: -----

PROGRESS: Number of examples : 17384

PROGRESS: Number of features : 6

PROGRESS: Number of unpacked features : 6

PROGRESS: Number of coefficients : 7

PROGRESS: Starting Newton Method

PROGRESS: -----

PROGRESS: +-----+-----+-----+-----+-----+-----+-----+
-----+

PROGRESS:	Iteration		Passes		Elapsed Time		Training-max_error	
Training-rmse								

PROGRESS: +-----+-----+-----+-----+-----+-----+-----+
-----+

PROGRESS:	1		2		0.019001		4014170.932927	
235190.935428								

PROGRESS: +-----+-----+-----+-----+-----+-----+-----+
-----+

PROGRESS: SUCCESS: Optimal solution found.

PROGRESS:

PROGRESS: Linear regression:

PROGRESS: -----

PROGRESS: Number of examples : 17384

PROGRESS: Number of features : 9

PROGRESS: Number of unpacked features : 9

PROGRESS: Number of coefficients : 10

PROGRESS: Starting Newton Method

PROGRESS: -----

PROGRESS: +-----+-----+-----+-----+-----+-----+-----+
-----+

```

PROGRESS: | Iteration | Passes | Elapsed Time | Training-max_error |
Training-rmse |
PROGRESS: +-----+-----+-----+-----+-----+
-----+
PROGRESS: | 1 | 2 | 0.009000 | 3193229.177894 |
228200.043155 |
PROGRESS: +-----+-----+-----+-----+-----+
-----+
PROGRESS: SUCCESS: Optimal solution found.
PROGRESS:

```

```

In [86]:
# Examine/extract each model's coefficients:
print model_1.get("coefficients")
print model_2.get("coefficients")
print model_3.get("coefficients")

```

name	index	value
(intercept)	None	-56140675.7444
sqft_living	None	310.263325778
bedrooms	None	-59577.1160682
bathrooms	None	13811.8405418
lat	None	629865.789485
long	None	-214790.285186

```
[6 rows x 3 columns]
```

name	index	value
(intercept)	None	-54410676.1152
sqft_living	None	304.449298057
bedrooms	None	-116366.043231
bathrooms	None	-77972.3305135
lat	None	625433.834953
long	None	-203958.60296
bed_bath_rooms	None	26961.6249092

```
[7 rows x 3 columns]
```

name	index	value
(intercept)	None	-52974974.0602
sqft_living	None	529.196420564
bedrooms	None	28948.5277313
bathrooms	None	65661.207231
lat	None	704762.148408
long	None	-137780.01994
bed_bath_rooms	None	-8478.36410518
bedrooms_squared	None	-6072.38466067
log_sqft_living	None	-563467.784269
lat_plus_long	None	-83217.1979248

```
+-----+-----+-----+
[10 rows x 3 columns]
```

Quiz Question: What is the sign (positive or negative) for the coefficient/weight for 'bathrooms' in model 1?

Quiz Question: What is the sign (positive or negative) for the coefficient/weight for 'bathrooms' in model 2?

Think about what this means.

Comparing multiple models

Now that you've learned three models and extracted the model weights we want to evaluate which model is best.

First use your functions from earlier to compute the RSS on TRAINING Data for each of the three models.

```
In [87]:
```

```
# Compute the RSS on TRAINING data for each of the three models and record the values:
```

```
print get_residual_sum_of_squares(model_1, train_data,
train_data['price'])
```

```
print get_residual_sum_of_squares(model_2, train_data,
train_data['price'])
```

```
print get_residual_sum_of_squares(model_3, train_data,
train_data['price'])
```

```
9.71328233544e+14
```

```
9.61592067856e+14
```

```
9.05276314555e+14
```

Quiz Question: Which model (1, 2 or 3) has lowest RSS on TRAINING Data? Is this what you expected?

Now compute the RSS on on TEST data for each of the three models.

```
In [88]:
```

```
# Compute the RSS on TESTING data for each of the three models and record the values:
```

```
print get_residual_sum_of_squares(model_1, test_data, test_data['price'])
```

```
print get_residual_sum_of_squares(model_2, test_data, test_data['price'])
```

```
print get_residual_sum_of_squares(model_3, test_data, test_data['price'])
```

```
2.26568089093e+14
```

```
2.24368799994e+14
```

```
2.51829318952e+14
```

Quiz Question: Which model (1, 2 or 3) has lowest RSS on TESTING Data? Is this what you expected? Think about the features that were added to each model from the previous.