

Fire up graphlab create

In [1]:

```
import graphlab
graphlab.canvas.set_target('ipynb')
```

Load some house sales data

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

In [2]:

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

```
/home/aluno/.virtualenvs/gl-env/local/lib/python2.7/site-packages/re
quests/packages/urllib3/util/ssl_.py:315: SNIMissingWarning: An HTTP
S request has been made, but the SNI (Subject Name Indication) exten
sion to TLS is not available on this platform. This may cause the se
rver to present an incorrect TLS certificate, which can cause valida
tion failures. For more information, see https://urllib3.readthedoc
s.org/en/latest/security.html#snimissingwarning.
```

```
SNIMissingWarning
```

```
/home/aluno/.virtualenvs/gl-env/local/lib/python2.7/site-packages/re
quests/packages/urllib3/util/ssl_.py:120: InsecurePlatformWarning: A
true SSLContext object is not available. This prevents urllib3 from
configuring SSL appropriately and may cause certain SSL connections
to fail. For more information, see https://urllib3.readthedocs.org/
en/latest/security.html#insecureplatformwarning.
```

```
InsecurePlatformWarning
```

This non-commercial license of GraphLab Create for academic use is assigned to alainandre@decom.cefetmg.br and will expire on July 19, 2017.

```
[INFO] graphlab.cython.cy_server: GraphLab Create v2.0.1 started. Lo
gging: /tmp/graphlab_server_1490113926.log
```

In [3]:

```
sales
```

Out[3]:

	00:00:00+00:00				
6414100192	2014-12-09 00:00:00+00:00	538000	3	2.25	2
5631500400	2015-02-25 00:00:00+00:00	180000	2	1	7
2487200875	2014-12-09 00:00:00+00:00	604000	4	3	1
1954400510	2015-02-18 00:00:00+00:00	510000	3	2	1
7237550310	2014-05-12 00:00:00+00:00	1225000	4	4.5	5
1321400060	2014-06-27 00:00:00+00:00	257500	3	2.25	1
2008000270	2015-01-15 00:00:00+00:00	291850	3	1.5	1
2414600126	2015-04-15 00:00:00+00:00	229500	3	1	1
3793500160	2015-03-12 00:00:00+00:00	323000	3	2.5	1

view	condition	grade	sqft_above	sqft_basement	yr_built	yr
0	3	7	1180	0	1955	
0	3	7	2170	400	1951	
0	3	6	770	0	1933	
0	5	7	1050	910	1965	
0	3	8	1680	0	1987	
0	3	11	3890	1530	2001	
0	3	7	1715	0	1995	
0	3	7	1060	0	1963	
0	3	7	1050	730	1960	
0	3	7	1890	0	2003	

long	sqft_living15	sqft_lot15
-122.25677536	1340.0	5650.0
-122.3188624	1690.0	7639.0
-122.23319601	2720.0	8062.0
-122.39318505	1360.0	5000.0

Exploring the data for housing sales

The house price is correlated with the number of square feet of living space.

In [6]:

```
sales.show(view="Scatter Plot", x="sqft_living", y="price")
```

Create a simple regression model of sqft_living to price

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 [7]:

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

Build the regression model using only sqft_living as a feature

In [8]:

```
sqft_model = graphlab.linear_regression.create(train_data,
                                              target='price',
                                              features=['sqft_living'],
                                              validation_set=None)
```

Linear regression:

```
-----

Number of examples      : 17384
Number of features      : 1
Number of unpacked features : 1
Number of coefficients   : 2
Starting Newton Method

-----

+-----+-----+-----+-----+-----+
-----+

| Iteration | Passes   | Elapsed Time | Training-max_error | Training-
rmse |
+-----+-----+-----+-----+-----+
-----+

| 1         | 2         | 1.013320     | 4349521.926170     | 262943.
613754 |
+-----+-----+-----+-----+-----+
-----+

SUCCESS: Optimal solution found.
```

Evaluate the simple model

In [9]:

```
print test_data['price'].mean()
```

543054.042563

In [10]:

```
print sqft_model.evaluate(test_data)
```

```
{'max_error': 4143550.8825285966, 'rmse': 255191.0287052736}
```

RMSE of about \$255,170!

Let's show what our predictions look like

Matplotlib is a Python plotting library that is also useful for plotting. You can install it with:

```
'pip install matplotlib'
```

In [13]:

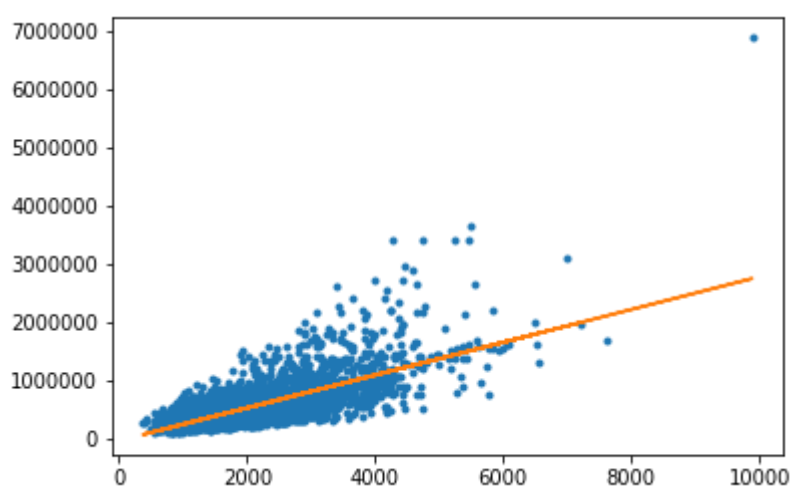
```
import matplotlib.pyplot as plt
%matplotlib inline
```

In [14]:

```
plt.plot(test_data['sqft_living'],test_data['price'],'.',
         test_data['sqft_living'],sqft_model.predict(test_data),'-')
```

Out[14]:

```
[<matplotlib.lines.Line2D at 0x7efdc40905d0>,
 <matplotlib.lines.Line2D at 0x7efdc4090750>]
```



Above: blue dots are original data, green line is the prediction from the simple regression.

Below: we can view the learned regression coefficients.

In [15]:

```
sqft_model.get('coefficients')
```

Out[15]:

name	index	value	stderr
(intercept)	None	-47114.0206702	4923.34437753
sqft_living	None	281.957850166	2.16405465323

[2 rows x 4 columns]

Explore other features in the data

To build a more elaborate model, we will explore using more features.

In [16]:

```
my_features = ['bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors', 'zipcode']
```

In [17]:

```
sales[my_features].show()
```

In [18]:

```
sales.show(view='BoxWhisker Plot', x='zipcode', y='price')
```

Pull the bar at the bottom to view more of the data.

98039 is the most expensive zip code.

Build a regression model with more features

In [19]:

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

Linear regression:

Number of examples : 17384

Number of features : 6

Number of unpacked features : 6

Number of coefficients : 115

Starting Newton Method

```
+-----+-----+-----+-----+-----+
+-----+
| Iteration | Passes   | Elapsed Time | Training-max_error | Trainin
g-rmse |
+-----+-----+-----+-----+-----+
+-----+
| 1         | 2        | 0.037419     | 3763208.270523     | 181908.
848367 |
+-----+-----+-----+-----+-----+
+-----+
```

SUCCESS: Optimal solution found.

In [20]:

```
print my_features
```

```
['bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors', 'zipc
ode']
```

Comparing the results of the simple model with adding more features

In [21]:

```
print sqft_model.evaluate(test_data)
print my_features_model.evaluate(test_data)
```

```
{'max_error': 4143550.8825285966, 'rmse': 255191.0287052736}
{'max_error': 3486584.509381705, 'rmse': 179542.4333126903}
```


The RMSE goes down from \$255,170 to \$179,508 with more features.

Apply learned models to predict prices of 3 houses

The first house we will use is considered an "average" house in Seattle.

In [22]:

```
house1 = sales[sales['id']=='5309101200']
```

In [23]:

```
house1
```

Out[23]:

id	date	price	bedrooms	bathrooms	sqft_liv
5309101200	2014-06-05 00:00:00+00:00	620000	4	2.25	2400

view	condition	grade	sqft_above	sqft_basement	yr_built	yr_renovated
0	4	7	1460	940	1929	

long	sqft_living15	sqft_lot15
-122.37010126	1250.0	4880.0

[? rows x 21 columns]

Note: Only the head of the SFrame is printed. This SFrame is lazily evaluated.

You can use `sf.materialize()` to force materialization.





In [24]:

```
print house1['price']  
[620000, ... ]
```

In [25]:

```
print sqft_model.predict(house1)  
[629584.8197281543]
```

In [26]:

```
print my_features_model.predict(house1)  
[721918.9333272863]
```

In this case, the model with more features provides a worse prediction than the simpler model with only 1 feature. However, on average, the model with more features is better.

Prediction for a second, fancier house

We will now examine the predictions for a fancier house.

In [27]:

```
house2 = sales[sales['id']=='1925069082']
```

In [28]:

```
house2
```

Out[28]:

id	date	price	bedrooms	bathrooms	sqft_li
1925069082	2015-05-11 00:00:00+00:00	2200000	5	4.25	464

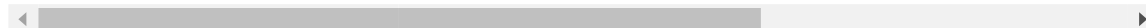
view	condition	grade	sqft_above	sqft_basement	yr_built	yr_re
4	5	8	2860	1780	1952	

long	sqft_living15	sqft_lot15
-122.09722322	3140.0	14200.0

[? rows x 21 columns]

Note: Only the head of the SFrame is printed. This SFrame is lazily evaluated.

You can use `sf.materialize()` to force materialization.





In [29]:

```
print house2['price']  
[2200000, ... ]
```

In [30]:

```
print sqft_model.predict(house2)  
[1261170.404099967]
```

In [31]:

```
print my_features_model.predict(house2)  
[1446472.4690774973]
```

In this case, the model with more features provides a better prediction. This behavior is expected here, because this house is more differentiated by features that go beyond its square feet of living space, especially the fact that it's a waterfront house.

Last house, super fancy

Our last house is a very large one owned by a famous Seattleite.

In [32]:

```
bill_gates = {'bedrooms':[8],  
              'bathrooms':[25],  
              'sqft_living':[50000],  
              'sqft_lot':[225000],  
              'floors':[4],  
              'zipcode':['98039'],  
              'condition':[10],  
              'grade':[10],  
              'waterfront':[1],  
              'view':[4],  
              'sqft_above':[37500],  
              'sqft_basement':[12500],  
              'yr_built':[1994],  
              'yr_renovated':[2010],  
              'lat':[47.627606],  
              'long':[-122.242054],  
              'sqft_living15':[5000],  
              'sqft_lot15':[40000]}
```



In [33]:

```
print my_features_model.predict(graphlab.SFrame(bill_gates))
```

```
[13749825.525719076]
```

The model predicts a price of over \$13M for this house! But we expect the house to cost much more. (There are very few samples in the dataset of houses that are this fancy, so we don't expect the model to capture a perfect prediction here.)

In [34]:

```
print sqft_model.predict(graphlab.SFrame(bill_gates))
```

```
[14050778.487629175]
```

Answers

1 - Selection and summary statistics

In [75]:

```
sales_filtered = [x for x in sales if x['zipcode'] == '98004']
```

In [94]:

```
avg = reduce(lambda x, y: x + y, [x['price'] for x in sales_filtered]) / len(sales_filtered)
avg
```

Out[94]:

1355927

2 - Filtering data

In [87]:

```
filtered_data = [x for x in sales if x['sqft_living'] >= 2000 and x['sqft_living'] <= 4000]
```

In [86]:

```
(1.0 * len(filtered_data)) / len(sales)
```

Out[86]:

0.4266413732475825

3 - Building a regression model with several more features

```
advanced_features = [
    'bedrooms',
    'bathrooms',
    'sqft_living',
    'sqft_lot',
    'floors',
    'zipcode',
    'condition', # condition of house
    'grade', # measure of quality of construction
    'waterfront', # waterfront property
    'view', # type of view
    'sqft_above', # square feet above ground
    'sqft_basement', # square feet in basement
    'yr_built', # the year built
    'yr_renovated', # the year renovated
    'lat', 'long', # the lat-long of the parcel
    'sqft_living15', # average sq.ft. of 15 nearest neighbors
    'sqft_lot15', # average lot size of 15 nearest neighbors
]
```

```
advanced_features_model = graphlab.linear_regression.create(train_data,target='price',features=advanced_features,validation set=None)
```

```

-----
Number of examples          : 17384
Number of features          : 18
Number of unpacked features : 18
Number of coefficients      : 127
Starting Newton Method
-----
+-----+-----+-----+-----+-----+
-----+
| Iteration | Passes   | Elapsed Time | Training-max_error | Training-
g-rmse |
+-----+-----+-----+-----+-----+
-----+
| 1         | 2        | 0.084595     | 3469012.450686    | 154580.
940736 |
+-----+-----+-----+-----+-----+
-----+
SUCCESS: Optimal solution found.

```