

Regression Week 3: Assessing Fit (polynomial regression)

In this notebook you will compare different regression models in order to assess which model fits best. We will be using polynomial regression as a means to examine this topic. In particular you will:

- Write a function to take an SArray and a degree and return an SFrame where each column is the SArray to a polynomial value up to the total degree e.g. degree = 3 then column 1 is the SArray column 2 is the SArray squared and column 3 is the SArray cubed
- Use matplotlib to visualize polynomial regressions
- Use matplotlib to visualize the same polynomial degree on different subsets of the data
- Use a validation set to select a polynomial degree
- Assess the final fit using test data

We will continue to use the House data from previous notebooks.

Fire up graphlab create

In [1]:

```
import graphlab
```

Next we're going to write a polynomial function that takes an SArray and a maximal degree and returns an SFrame with columns containing the SArray to all the powers up to the maximal degree.

The easiest way to apply a power to an SArray is to use the `.apply()` and `lambda x:` functions. For example to take the example array and compute the third power we can do as follows: (note running this cell the first time may take longer than expected since it loads graphlab)

In [2]:

```
tmp = graphlab.SArray([1., 2., 3.])
tmp_cubed = tmp.apply(lambda x: x**3)
print tmp
print tmp_cubed
```

```
[INFO] 1449568855 : INFO:      (initialize_globals_from_environment:28
2): Setting configuration variable GRAPHLAB_FILEIO_ALTERNATIVE_SSL_CER
T_FILE to C:\Users\Sandipan.Dey\AppData\Local\Dato\Dato Launcher\lib\si
te-packages\certifi\cacert.pem
1449568855 : INFO:      (initialize_globals_from_environment:282): Setti
ng configuration variable GRAPHLAB_FILEIO_ALTERNATIVE_SSL_CERT_DIR to
This non-commercial license of GraphLab Create is assigned to sandipa
n.dey@gmail.com and will expire on October 12, 2016. For commercial lic
ensing options, visit https://dato.com/buy/. (https://dato.com/buy/.)
```

```
[INFO] Start server at: ipc:///tmp/graphlab_server-11328 - Server binar
y: C:\Users\Sandipan.Dey\AppData\Local\Dato\Dato Launcher\lib\site-pack
ages\graphlab\unity_server.exe - Server log: C:\Users\Sandipan.Dey\AppData\Local\Temp\graphlab_server_1449568855.log.0
[INFO] GraphLab Server Version: 1.7.1
```

```
[1.0, 2.0, 3.0]
[1.0, 8.0, 27.0]
```

We can create an empty SFrame using `graphlab.SFrame()` and then add any columns to it with `ex_sframe['column_name'] = value`. For example we create an empty SFrame and make the column 'power_1' to be the first power of tmp (i.e. tmp itself).

In [3]:

```
ex_sframe = graphlab.SFrame()
ex_sframe['power_1'] = tmp
print ex_sframe
```

```
+-----+
| power_1 |
+-----+
|   1.0   |
|   2.0   |
|   3.0   |
+-----+
[3 rows x 1 columns]
```

Polynomial_sframe function

Using the hints above complete the following function to create an SFrame consisting of the powers of an SArray up to a specific degree:

In [8]:

```
def polynomial_sframe(feature, degree):
    # assume that degree >= 1
    # initialize the SFrame:
    poly_sframe = graphlab.SFrame()
    # and set poly_sframe['power_1'] equal to the passed feature
    poly_sframe['power_1'] = feature
    # first check if degree > 1
    if degree > 1:
        # then loop over the remaining degrees:
        # range usually starts at 0 and stops at the endpoint-1. We want it to start
        for power in range(2, degree+1):
            # first we'll give the column a name:
            name = 'power_' + str(power)
            # then assign poly_sframe[name] to the appropriate power of feature
            poly_sframe[name] = feature.apply(lambda x: x**power)
    return poly_sframe
```

To test your function consider the smaller tmp variable and what you would expect the outcome of the following call:

In [9]:

```
print polynomial_sframe(tmp, 3)
```

```
+-----+-----+-----+
| power_1 | power_2 | power_3 |
+-----+-----+-----+
|    1.0   |    1.0   |    1.0   |
|    2.0   |    4.0   |    8.0   |
|    3.0   |    9.0   |   27.0   |
+-----+-----+-----+
[3 rows x 3 columns]
```

Visualizing polynomial regression

Let's use matplotlib to visualize what a polynomial regression looks like on some real data.

In [10]:

```
sales = graphlab.SFrame('C:/courses/Coursera/Current/ML Regression/Week3/kc_house_dat
```

As in Week 3, we will use the `sqft_living` variable. For plotting purposes (connecting the dots), you'll need to sort by the values of `sqft_living`. For houses with identical square footage, we break the tie by their prices.

In [14]:

```
#Let's take a look at the weights before we plot
model1.get("coefficients")
```

Out[14]:

name	index	value
(intercept)	None	-43579.0852515
power_1	None	280.622770886

[2 rows x 3 columns]

In [15]:

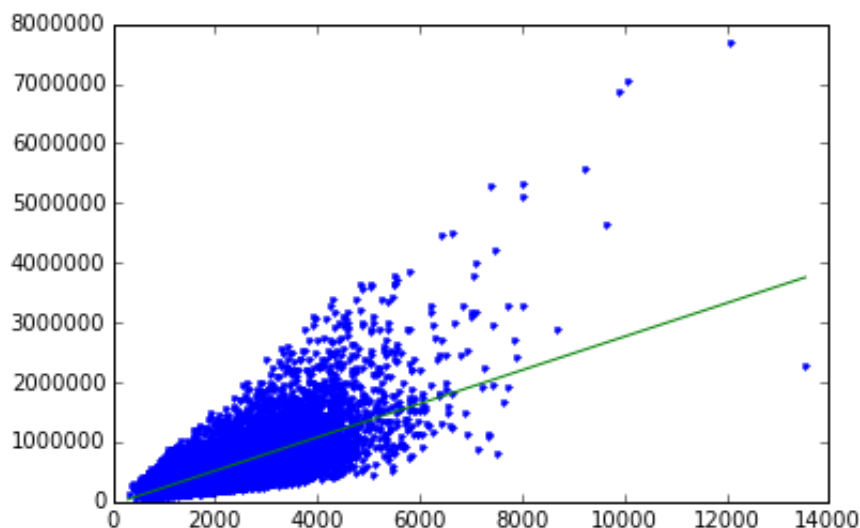
```
import matplotlib.pyplot as plt
%matplotlib inline
```

In [16]:

```
plt.plot(poly1_data['power_1'],poly1_data['price'],'.',
         poly1_data['power_1'], model1.predict(poly1_data),'-')
```

Out[16]:

```
[<matplotlib.lines.Line2D at 0x1ee30b70>,
 <matplotlib.lines.Line2D at 0x1ee30da0>]
```

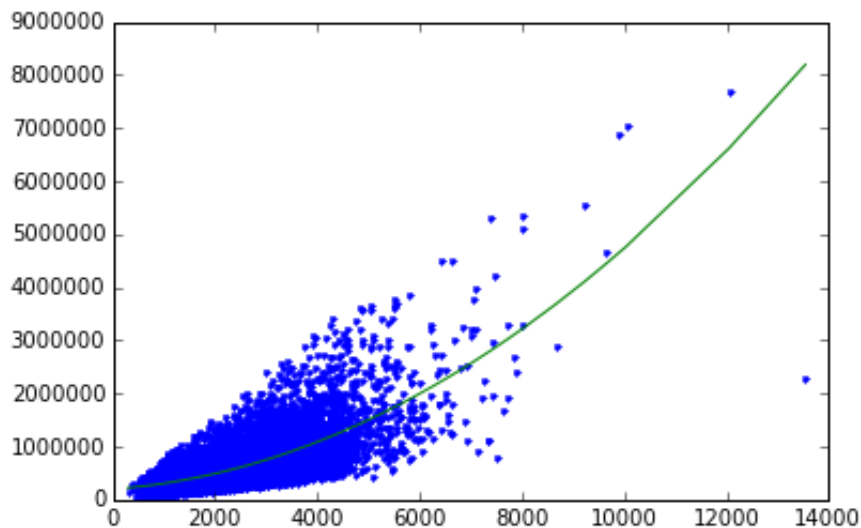


Let's unpack that `plt.plot()` command. The first pair of SArrays we passed are the 1st power of sqft and the actual price we then ask it to print these as dots `'.'`. The next pair we pass is the 1st power of sqft and the predicted values from the linear model. We ask these to be plotted as a line `'-'`.

We can see, not surprisingly, that the predicted values all fall on a line, specifically the one with slope 280 and intercept -43579. What if we wanted to plot a second degree polynomial?


```
plt.plot(poly2_data['power_1'],poly2_data['price'],'.',
        poly2_data['power_1'], model2.predict(poly2_data),'-')
```

```
[<matplotlib.lines.Line2D at 0x1f2936d8>,  
 <matplotlib.lines.Line2D at 0x1f2938d0>]
```



```
poly2_data = polynomial_sframe(sales['sqft_living'], 3)
my_features = poly2_data.column_names() # get the name of the features
poly2_data['price'] = sales['price'] # add price to the data since it's the target
model2 = graphlab.linear_regression.create(poly2_data, target = 'price', features = m
```

```

PROGRESS: -----
PROGRESS: Number of examples          : 21613
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          | 0.000000    | 3261066.736007    |
249261.286346 |
PROGRESS: +-----+-----+-----+-----+-----+
-----+
PROGRESS: SUCCESS: Optimal solution found.
PROGRESS:

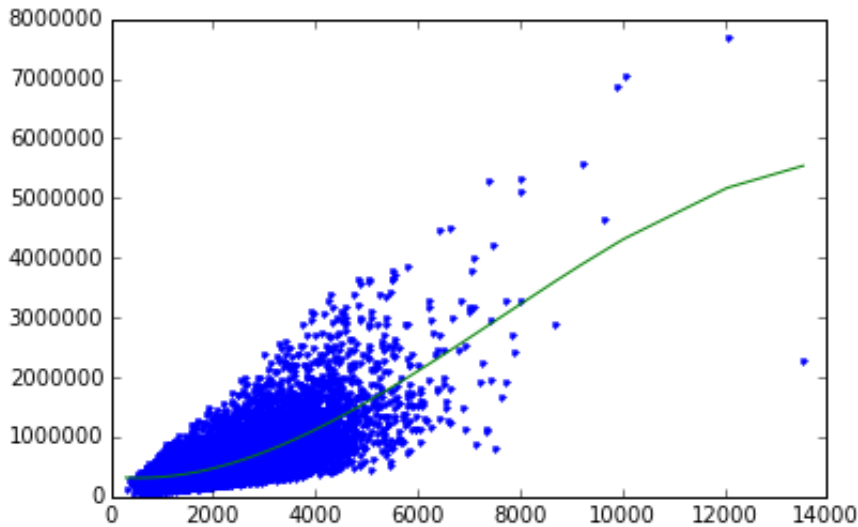
```

In [21]:

```
plt.plot(poly2_data['power_1'],poly2_data['price'],'.',
         poly2_data['power_1'], model2.predict(poly2_data),'-')
```

Out[21]:

```
[<matplotlib.lines.Line2D at 0x1f42ffd0>,
 <matplotlib.lines.Line2D at 0x1f43f208>]
```



Now try a 15th degree polynomial:

In [22]:

```
poly2_data = polynomial_sframe(sales['sqft_living'], 15)
my_features = poly2_data.column_names() # get the name of the features
poly2_data['price'] = sales['price'] # add price to the data since it's the target
model2 = graphlab.linear_regression.create(poly2_data, target = 'price', features = m
```

PROGRESS: Linear regression:

PROGRESS: -----

PROGRESS: Number of examples : 21613

PROGRESS: Number of features : 15

PROGRESS: Number of unpacked features : 15

PROGRESS: Number of coefficients : 16

PROGRESS: Starting Newton Method

PROGRESS: -----

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

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

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

```
PROGRESS: | 1 | 2 | 0.012728 | 2662308.584338 |
245690.511190 |
```

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

PROGRESS: SUCCESS: Optimal solution found.

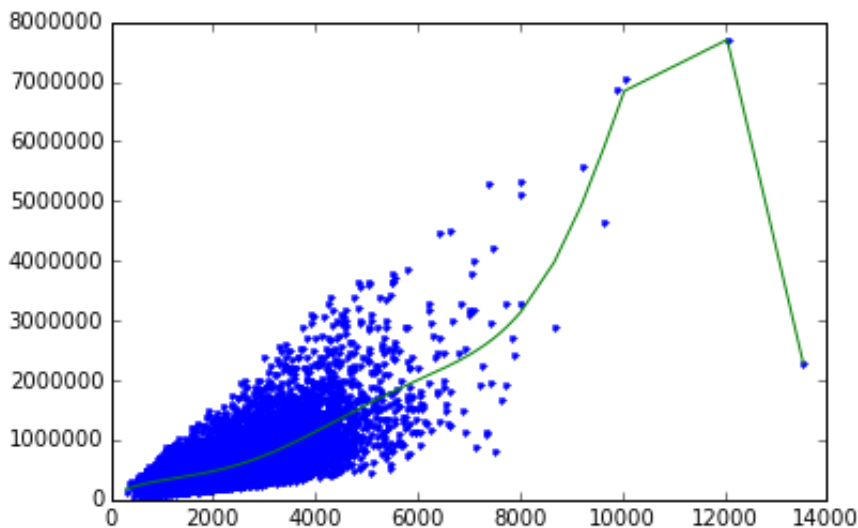
PROGRESS:

In [23]:

```
plt.plot(poly2_data['power_1'],poly2_data['price'],'.',
         poly2_data['power_1'], model2.predict(poly2_data),'-')
```

Out[23]:

```
[<matplotlib.lines.Line2D at 0x1f9bce80>,
 <matplotlib.lines.Line2D at 0x1f9cc0b8>]
```



What do you think of the 15th degree polynomial? Do you think this is appropriate? If we were to change the data do you think you'd get pretty much the same curve? Let's take a look.

Changing the data and re-learning

We're going to split the sales data into four subsets of roughly equal size. Then you will estimate a 15th degree polynomial model on all four subsets of the data. Print the coefficients (you should use `.print_rows(num_rows = 16)` to view all of them) and plot the resulting fit (as we did above). The quiz will ask you some questions about these results.

To split the sales data into four subsets, we perform the following steps:

- First split sales into 2 subsets with `.random_split(0.5, seed=0)`.
- Next split the resulting subsets into 2 more subsets each. Use `.random_split(0.5, seed=0)`.

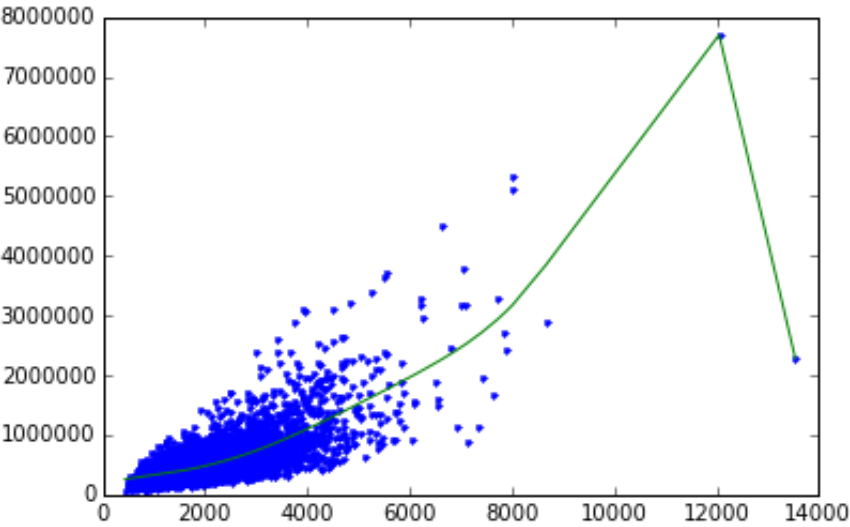
We set `seed=0` in these steps so that different users get consistent results. You should end up with 4 subsets (`set_1`, `set_2`, `set_3`, `set_4`) of approximately equal size.

In [24]:

```
data_1,data_2 = sales.random_split(.5,seed=0)
set_1,set_2 = data_1.random_split(.5,seed=0)
set_3,set_4 = data_2.random_split(.5,seed=0)
```

Fit a 15th degree polynomial on `set_1`, `set_2`, `set_3`, and `set_4` using `sqft_living` to predict prices. Print

the coefficients and make a plot of the resulting model.



```
poly_data = polynomial_sframe(set_2['sqft_living'], 15)
my_features = poly_data.column_names() # get the name of the features
poly_data['price'] = set_2['price'] # add price to the data since it's the target
model_2 = graphlab.linear_regression.create(poly_data, target = 'price', features = my_features)
model_2.get("coefficients").print_rows(num_rows=16)
plt.plot(poly_data['power_1'], poly_data['price'], '.', poly_data['power_1'], model_2.predict(poly_data['power_1']))
```

PROGRESS: -----

PROGRESS: Number of features : 15

PROGRESS: Number of coefficients : 16

PROGRESS: -----

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

-----+

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

-----+

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

-----+

PROGRESS:

+

name	index	value
------	-------	-------

+

(intercept)	None	160515.194669
-------------	------	---------------

power 1	None	161.068906213
---------	------	---------------

power 2	None	0.00721288554281
---------	------	------------------

power 3	None	-1.53767451328e-05
---------	------	--------------------

power 4	None	5.53101276894e-09
---------	------	-------------------

power 5	None	3.44914141903e-13
---------	------	-------------------

power 6	None	-8.41349331246e-17
---------	------	--------------------

power 7	None	-1.17557544119e-20
---------	------	--------------------

power 8	None	-3.24855695819e-25
---------	------	--------------------

power 9	None	8.06950508874e-29
---------	------	-------------------

power_10	None	1.36060382528e-32
----------	------	-------------------

power_11	None	1.06720789566e-36
----------	------	-------------------

power_12	None	1.92370157332e-41
----------	------	-------------------

power_13	None	-7.10368037472e-45
----------	------	--------------------

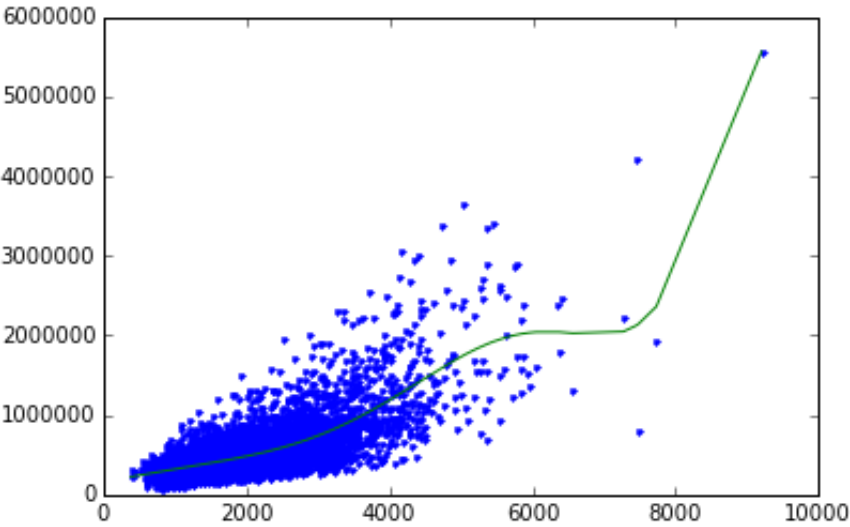
power_14	None	-1.01868434954e-48
----------	------	--------------------

power_15	None	-1.74692416518e-53
----------	------	--------------------

```
[16 rows x 3 columns]
```

```
[<matplotlib.lines.Line2D at 0x20ed3940>,
```

```
<matplotlib.lines.Line2D at 0x20ed3b38>]
```



```
poly_data = polynomial_sframe(set_3['sqft_living'], 15)
my_features = poly_data.column_names() # get the name of the features
poly_data['price'] = set_3['price'] # add price to the data since it's the target
model_3 = graphlab.linear_regression.create(poly_data, target = 'price', features = my_features)
model_3.get("coefficients").print_rows(num_rows=16)
plt.plot(poly_data['power_1'], poly_data['price'], '.', poly_data['power_1'], model_3.predict(poly_data['power_1']))
```

PROGRESS: -----

PROGRESS: Number of features : 15

PROGRESS: Number of coefficients : 16

PROGRESS: -----

-----+

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

-----+

247777.375090 |

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

-----+

PROGRESS:

-----+

name	index	value
------	-------	-------

+

(intercept)	None	64031.5743611
-------------	------	---------------

power 1	None	419.963446533
---------	------	---------------

power 2	None	-0.217032383683
---------	------	-----------------

power 3	None	5.71721871042e-05
---------	------	-------------------

power 4	None	6.42456678907e-10
---------	------	-------------------

power 5	None	-8.76764336026e-13
---------	------	--------------------

power 6	None	-4.39079425e-17
---------	------	-----------------

power 7	None	4.65780734822e-21
---------	------	-------------------

power 8	None	8.56537636021e-25
---------	------	-------------------

power 9	None	5.87177076692e-29
---------	------	-------------------

power_10	None	-1.18822533167e-34
----------	------	--------------------

power_11	None	-5.48520101257e-37
----------	------	--------------------

power_12	None	-7.6344609578e-41
----------	------	-------------------

power_13	None	-5.62621010165e-45
----------	------	--------------------

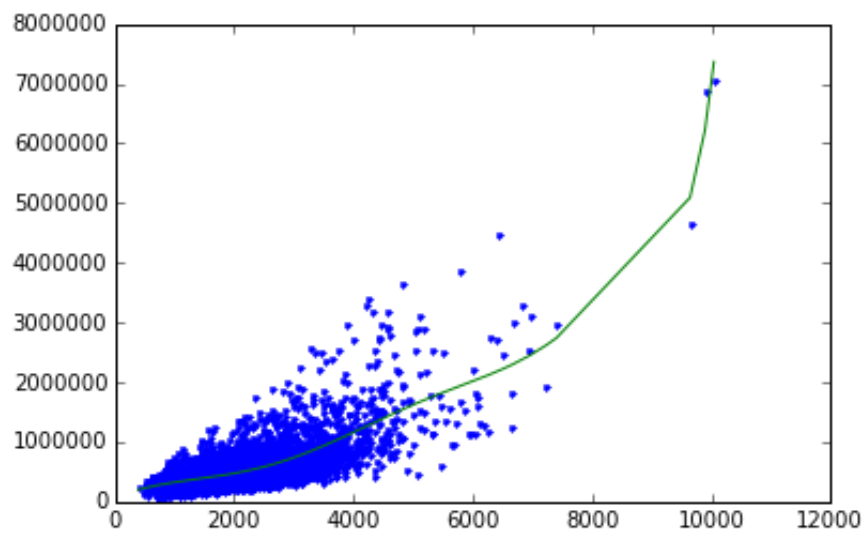
power_14	None	2.10087982472e-50
----------	------	-------------------

power_15	None	8.8038750037e-53
----------	------	------------------

```
[16 rows x 3 columns]
```

```
[<matplotlib.lines.Line2D at 0x21593588>,
```

```
<matplotlib.lines.Line2D at 0x21593780>]
```




```
poly_data = polynomial_sframe(set_4['sqft_living'], 15)
my_features = poly_data.column_names() # get the name of the features
poly_data['price'] = set_4['price'] # add price to the data since it's the target
model_4 = graphlab.linear_regression.create(poly_data, target = 'price', features = my_features)
model_4.get("coefficients").print_rows(num_rows=16)
plt.plot(poly_data['power_1'], poly_data['price'], '.', poly_data['power_1'], model_4.predict(poly_data['power_1']))
```

PROGRESS: -----

PROGRESS: -----

-----+

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

-----+

242204.386550 |

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

-----+

PROGRESS:

-----+

name	index	value
------	-------	-------

-----+

(intercept)	None	238215.539489
-------------	------	---------------

power 1	None	35.6890462022
---------	------	---------------

power 2	None	0.0384180337882
---------	------	-----------------

power 3	None	1.00407290036e-05
---------	------	-------------------

power 4	None	-5.35136998633e-09
---------	------	--------------------

power 5	None	3.35662136037e-13
---------	------	-------------------

power 6	None	1.81755721542e-16
---------	------	-------------------

power 7	None	6.62015234312e-21
---------	------	-------------------

power 8	None	-3.13250645184e-24
---------	------	--------------------

power 9	None	-6.11495895413e-28
---------	------	--------------------

power_10	None	-4.37319305466e-32
----------	------	--------------------

power_11	None	3.5466600452e-36
----------	------	------------------

power_12	None	1.47563558469e-39
----------	------	-------------------

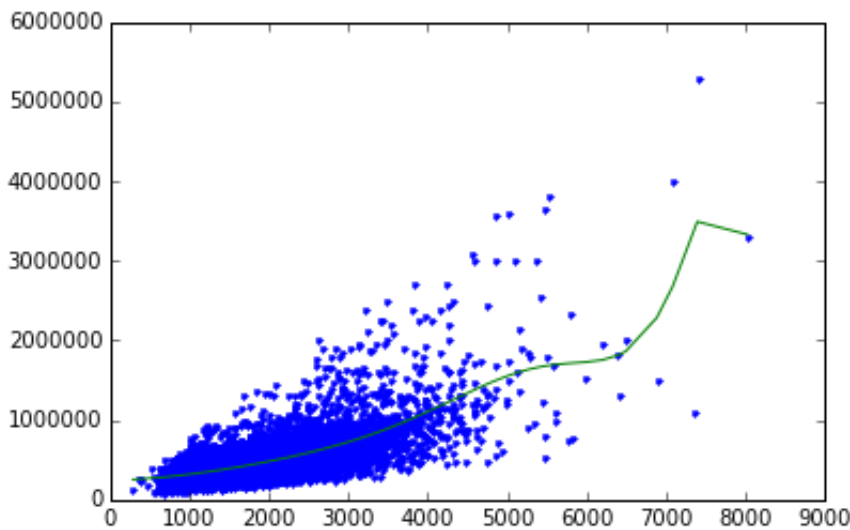
power_13	None	2.10094882436e-43
----------	------	-------------------

power_14	None	8.12201995016e-48
----------	------	-------------------

power_15	None	-4.45547001483e-51
----------	------	--------------------

```
[16 rows x 3 columns]
```

```
<matplotlib.lines.Line2D at 0x217bbcf8>]
```



Some questions you will be asked on your quiz:

Quiz Question: Is the sign (positive or negative) for `power_15` the same in all four models?

Quiz Question: (True/False) the plotted fitted lines look the same in all four plots

Selecting a Polynomial Degree

Whenever we have a "magic" parameter like the degree of the polynomial there is one well-known way to select these parameters: validation set. (We will explore another approach in week 4).

We split the sales dataset 3-way into training set, test set, and validation set as follows:

```
* Split our sales data into 2 sets: `training_and_validation` and `testing`. Use
`random_split(0.9, seed=1)`.
* Further split our training data into two sets: `training` and `validation`. Use
`random_split(0.5, seed=1)`.
```

Again, we set `seed=1` to obtain consistent results for different users.

In [41]:

```
training_and_validation,testing = sales.random_split(.9,seed=1)
training,validation = training_and_validation.random_split(.5,seed=1)
```

Next you should write a loop that does the following:

- For degree in [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15] (to get this in python type `range(1, 15+1)`)
 - Build an SFrame of polynomial data of `train_data['sqft_living']` at the current degree
 - hint: `my_features = poly_data.column_names()` gives you a list e.g. `['power_1', 'power_2', 'power_3']` which you might find useful for `graphlab.linear_regression.create(features = my_features)`
 - Add `train_data['price']` to the polynomial SFrame

- Learn a polynomial regression model to sqft vs price with that degree on TRAIN data
- Compute the RSS on VALIDATION data (here you will want to use `.predict()`) for that degree and you will need to make a polynomial SFrame using validation data.
- Report which degree had the lowest RSS on validation data (remember python indexes from 0)

(Note you can turn off the print out of `linear_regression.create()` with `verbose = False`)

In [51]:

```
RSSVals = {}
for degree in range(1, 16):
    poly_data = polynomial_sframe(training['sqft_living'], degree)
    my_features = poly_data.column_names() # get the name of the features
    poly_data['price'] = training['price'] # add price to the data since it's the target
    model = graphlab.linear_regression.create(poly_data, target = 'price', features =
    #model.get("coefficients").print_rows(num_rows=degree)
    poly_data = polynomial_sframe(validation['sqft_living'], degree)
    RSS=sum((validation['price'] - model.predict(poly_data))**2)
    print degree, RSS
    RSSVals[degree] = RSS
    #plt.plot(poly_data['power_1'],poly_data['price'],'.', poly_data['power_1'], model.predict(poly_data))
for degree in RSSVals:
    print degree, RSSVals[degree]
print min([(RSSVals[degree], degree) for degree in RSSVals])
```

```
--+-----+
PROGRESS: SUCCESS: Optimal solution found.
PROGRESS:
11 6.36152414427e+14
PROGRESS: Linear regression:
PROGRESS: -----
PROGRESS: Number of examples      : 9761
PROGRESS: Number of features      : 12
PROGRESS: Number of unpacked features : 12
PROGRESS: Number of coefficients   : 13
PROGRESS: Starting Newton Method
PROGRESS: -----
PROGRESS: +-----+-----+-----+-----+
--+-----+
PROGRESS: | Iteration | Passes   | Elapsed Time | Training-max_error
| Training-rmse |
PROGRESS: +-----+-----+-----+-----+
--+-----+
PROGRESS: | 1          | 2        | 0.020004    | 2538026.076279
| 213270 176881 |
```

Quiz Question: Which degree (1, 2, ..., 15) had the lowest RSS on Validation data?

Now that you have chosen the degree of your polynomial using validation data, compute the RSS of this model on TEST data. Report the RSS on your quiz.

```
degree = 6
poly_data = polynomial_sframe(training['sqft_living'], degree)
my_features = poly_data.column_names() # get the name of the features
poly_data['price'] = training['price'] # add price to the data since it's the target
model = graphlab.linear_regression.create(poly_data, target = 'price', features = my_
#model.get("coefficients").print_rows(num_rows=degree)
poly_data = polynomial_sframe(testing['sqft_living'], degree)
RSS=sum((testing['price'] - model.predict(poly_data))**2)
print RSS
```

```

PROGRESS: Linear regression:
PROGRESS: -----
PROGRESS: Number of examples      : 9761
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          | 0.010002     | 2359346.812645     |
243809.143530 |
PROGRESS: +-----+-----+-----+-----+
-----+
PROGRESS: SUCCESS: Optimal solution found.
PROGRESS:
1.28190059156e+14

```

In []: