

Getting Started

To familiarize yourself with an iPython Notebook, **try double clicking on this cell**. You will notice that the text changes so that all the formatting is removed. This allows you to make edits to the block of text you see here. This block of text (and mostly anything that's not code) is written using Markdown (<http://daringfireball.net/projects/markdown/syntax>), which is a way to format text using headers, links, italics, and many other options! Whether you're editing a Markdown text block or a code block (like the one below), you can use the keyboard shortcut **Shift + Enter** or **Shift + Return** to execute the code or text block. In this case, it will show the formatted text.

Let's start by setting up some code we will need to get the rest of the project up and running. Use the keyboard shortcut mentioned above on the following code block to execute it. Alternatively, depending on your iPython Notebook program, you can press the **Play** button in the hotbar. You'll know the code block executes successfully if the message *"Boston Housing dataset loaded successfully!"* is printed.

In [4]:

```
# Importing a few necessary libraries
import numpy as np
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn.tree import DecisionTreeRegressor

# Make matplotlib show our plots inline (nicely formatted in the notebook)
%matplotlib inline

# Create our client's feature set for which we will be predicting a selling price
CLIENT_FEATURES = [[11.95, 0.00, 18.100, 0, 0.6590, 5.6090, 90.00, 1.385, 24, 680.0, 20.20, 332.09, 12.13]]

# Load the Boston Housing dataset into the city_data variable
city_data = datasets.load_boston()

# Initialize the housing prices and housing features
#print city_data
housing_prices = city_data.target
housing_features = city_data.data
print "Boston Housing dataset loaded successfully!"
```

Boston Housing dataset loaded successfully!

Statistical Analysis and Data Exploration

In this first section of the project, you will quickly investigate a few basic statistics about the dataset you are working with. In addition, you'll look at the client's feature set in `CLIENT_FEATURES` and see how this particular sample relates to the features of the dataset. Familiarizing yourself with the data through an explorative process is a fundamental practice to help you better understand your results.

Step 1

In the code block below, use the imported numpy library to calculate the requested statistics. You will need to replace each None you find with the appropriate numpy coding for the proper statistic to be printed. Be sure to execute the code block each time to test if your implementation is working successfully. The print statements will show the statistics you calculate!

In [5]:

```
# Number of houses in the dataset
total_houses = np.size(housing_prices)

# Number of features in the dataset
total_features = housing_features.shape[1]

# Minimum housing value in the dataset
minimum_price = np.amin(housing_prices)

# Maximum housing value in the dataset
maximum_price = np.amax(housing_prices)

# Mean house value of the dataset
mean_price = np.mean(housing_prices)

# Median house value of the dataset
median_price = np.median(housing_prices)

# Standard deviation of housing values of the dataset
std_dev = np.std(housing_prices)

# Show the calculated statistics
print "Boston Housing dataset statistics (in $1000's):\n"
print "Total number of houses:", total_houses
print "Total number of features:", total_features
print "Minimum house price:", minimum_price
print "Maximum house price:", maximum_price
print "Mean house price: {0:.3f}".format(mean_price)
print "Median house price:", median_price
print "Standard deviation of house price: {0:.3f}".format(std_dev)
```

Boston Housing dataset statistics (in \$1000's):

```
Total number of houses: 506
Total number of features: 13
Minimum house price: 5.0
Maximum house price: 50.0
Mean house price: 22.533
Median house price: 21.2
Standard deviation of house price: 9.188
```

Question 1

As a reminder, you can view a description of the Boston Housing dataset [here](https://archive.ics.uci.edu/ml/datasets/Housing) (<https://archive.ics.uci.edu/ml/datasets/Housing>), where you can find the different features under **Attribute Information**. The MEDV attribute relates to the values stored in our `housing_prices` variable, so we do not consider that a feature of the data.

Of the features available for each data point, choose three that you feel are significant and give a brief description for each of what they measure.

Remember, you can **double click the text box below** to add your answer!

1. CRIM: per capita crime rate by town - Crime rate is one of the important factor that make housing price increase and decrease. No one wants to buy an expensive house where around their house is so dangerous.
2. RM: average number of rooms per dwelling - It's obvious that even search engine or search criteria use number of rooms to measure the price of the house but pricing will not increase dramatically when there are too many rooms in the house. I believe that 5 Bedroom, in my opinion, is the most biggest of normal house after that pricing won't increase that much
3. DIS: weighted distances to five Boston employment centres - You spend half of your life at home and another half of your life at work. People are focusing to live happily close to their office also. That's why the closer of the downtown, where they are a lot of jobs, the more expensive it is because everyone thinks that same, live close to work place.

Question 2

Using your client's feature set `CLIENT_FEATURES`, which values correspond with the features you've chosen above?

Hint: Run the code block below to see the client's data.

In [6]:

```
print CLIENT_FEATURES
```

```
[[11.95, 0.0, 18.1, 0, 0.659, 5.609, 90.0, 1.385, 24, 680.0, 2  
0.2, 332.09, 12.13]]
```

11.95 - CRIM 5.609 - RM 1.385 - DIS

Evaluating Model Performance

In this second section of the project, you will begin to develop the tools necessary for a model to make a prediction. Being able to accurately evaluate each model's performance through the use of these tools helps to greatly reinforce the confidence in your predictions.

Step 2

In the code block below, you will need to implement code so that the `shuffle_split_data` function does the following:

- Randomly shuffle the input data `x` and target labels (housing values) `y`.
- Split the data into training and testing subsets, holding 30% of the data for testing.

If you use any functions not already accessible from the imported libraries above, remember to include your import statement below as well!

Ensure that you have executed the code block once you are done. You'll know the `shuffle_split_data` function is working if the statement *"Successfully shuffled and split the data!"* is printed.

In [20]:

```
# Put any import statements you need for this code block here
from sklearn.cross_validation import train_test_split

def shuffle_split_data(X, y):
    """ Shuffles and splits data into 70% training and 30% testing subsets,
        then returns the training and testing subsets. """

    # Shuffle and split the data
    #X_train = None
    #y_train = None
    #X_test = None
    #y_test = None
    # use sklearn algorithm which should be faster than my logic
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.
3, random_state=42)
    # Return the training and testing data subsets
    return X_train, y_train, X_test, y_test

# Test shuffle_split_data
try:
    X_train, y_train, X_test, y_test = shuffle_split_data(housing_features,
housing_prices)
    print "Successfully shuffled and split the data!"
except:
    print "Something went wrong with shuffling and splitting the data."
```

Successfully shuffled and split the data!

Question 3

Why do we split the data into training and testing subsets for our model?

So that when we want to validate our accuracy we could use testing data to test it. No need to find new dataset to test

Step 3

In the code block below, you will need to implement code so that the `performance_metric` function does the following:

- Perform a total error calculation between the true values of the `y` labels `y_true` and the predicted values of the `y` labels `y_predict`.

You will need to first choose an appropriate performance metric for this problem. See [the sklearn metrics documentation \(http://scikit-learn.org/stable/modules/classes.html#sklearn-metrics-metrics\)](http://scikit-learn.org/stable/modules/classes.html#sklearn-metrics-metrics) to view a list of available metric functions. **Hint:** Look at the question below to see a list of the metrics that were covered in the supporting course for this project.

Once you have determined which metric you will use, remember to include the necessary import statement as well!

Ensure that you have executed the code block once you are done. You'll know the `performance_metric` function is working if the statement *"Successfully performed a metric calculation!"* is printed.

In [44]:

```
# Put any import statements you need for this code block here
import numpy as np

from sklearn.metrics import mean_squared_error
from sklearn.metrics import mean_absolute_error

def performance_metric(y_true, y_predict):
    """ Calculates and returns the total error between true and predicted v
    alues
        based on a performance metric chosen by the student. """
    error = mean_squared_error(y_true, y_predict)
    #error = mean_absolute_error(y_true, y_predict)
    #print error
    return error

# Test performance_metric
try:
    total_error = performance_metric(y_train, y_train)
    print "Successfully performed a metric calculation!"
except:
    print "Something went wrong with performing a metric calculation."
```

Successfully performed a metric calculation!

Question 4

Which performance metric below did you find was most appropriate for predicting housing prices and analyzing the total error. Why?

- Accuracy
- Precision
- Recall
- F1 Score
- Mean Squared Error (MSE)
- Mean Absolute Error (MAE)

boston housing count as regression because we have a data which is continuous so we're going to use regression for analyzing, which MSE and MAE. Between MSE and MAE, I pick because MSE will be focus on giving more weight to points where further away from the mean which is make sense in most cases for checking the error while MAE are considered giving less weight when compare to MSE

Step 4 (Final Step)

In the code block below, you will need to implement code so that the `fit_model` function does the following:

- Create a scoring function using the same performance metric as in **Step 3**. See the [sklearn make_scorer](http://scikit-learn.org/stable/modules/generated/sklearn.metrics.make_scorer.html) documentation (http://scikit-learn.org/stable/modules/generated/sklearn.metrics.make_scorer.html).
- Build a GridSearchCV object using `regressor`, `parameters`, and `scoring_function`. See the [sklearn documentation on GridSearchCV](http://scikit-learn.org/stable/modules/generated/sklearn.grid_search.GridSearchCV.html) (http://scikit-learn.org/stable/modules/generated/sklearn.grid_search.GridSearchCV.html).

When building the scoring function and GridSearchCV object, *be sure that you read the parameters documentation thoroughly*. It is not always the case that a default parameter for a function is the appropriate setting for the problem you are working on.

Since you are using `sklearn` functions, remember to include the necessary import statements below as well!

Ensure that you have executed the code block once you are done. You'll know the `fit_model` function is working if the statement *"Successfully fit a model to the data!"* is printed.

In [51]:

```
# Put any import statements you need for this code block
from sklearn.metrics import make_scorer
from sklearn.grid_search import GridSearchCV
from sklearn.metrics import mean_squared_error
from sklearn.cross_validation import KFold

def fit_model(X, y):
    """ Tunes a decision tree regressor model using GridSearchCV on the input data X
        and target labels y and returns this optimal model. """

    # Create a decision tree regressor object
    regressor = DecisionTreeRegressor()
    A = [1,2,3,4,5]
    # Set up the parameters we wish to tune
    parameters = {'max_depth':(1,2,3,4,5,6,7,8,9,10)}
    # Make an appropriate scoring function
    scoring_function = make_scorer(performance_metric, greater_is_better=False)

    # Make the GridSearchCV object
    reg = GridSearchCV(regressor, parameters, scoring=scoring_function)

    # Fit the learner to the data to obtain the optimal model with tuned parameters
    reg.fit(X, y)
    #print reg.best_estimator_
    # Return the optimal model
    return reg.best_estimator_

# Test fit_model on entire dataset
try:
    reg = fit_model(housing_features, housing_prices)
    print "Successfully fit a model!"
except:
    print "Something went wrong with fitting a model."
```

Successfully fit a model!

Question 5

What is the grid search algorithm and when is it applicable?

In []:

Grid search algorithm is the traditional way of performing hyperparameter optimization. Grid search algorithm is applicable when small number of design variable in the datasets and the dataset is not too complexity and suffers from the curse of dimensionality. As I've mentioned that, It's not suitable for dataset which not too complexity it consumes too much time and space to calculate such as XOR Model.

Question 6

What is cross-validation, and how is it performed on a model? Why would cross-validation be helpful when using grid search?

Cross-validation is the validation technique that make data independent from each others so user can estimate the accuracy and prediction of the model in the new set(testing set). Cross-validation divided dataset into 2 parts, training set and testing set. Data will be train in the training set while validate the accuracy in the testing set. Cross-validation will be helpful on all algorithms for testing the accuracy not only with grid search. There some technique in cross-validation that used to generate split data, such as K-fold which the concept is divided dataset into k group and use k-1 to be the training set data while the rest used for validate or testing data. For example we have 300 dataset. If we are using normal split data we will have 210 in training set(70%) and 90 in testing sets, and we could testing the data just only once. In the meantime, K-fold will divide data into k group, let's say 10 group. Each group will have 30 data then we will have 270 training set and 30 testing set in the first round. After that we split again with the second K-fold so we will have completely new 270 training set with 30 testing set again to test your accuracy. so in theory, you have k-fold times to test your training set with testing set.

Checkpoint!

You have now successfully completed your last code implementation section. Pat yourself on the back! All of your functions written above will be executed in the remaining sections below, and questions will be asked about various results for you to analyze. To prepare the **Analysis** and **Prediction** sections, you will need to initialize the two functions below. Remember, there's no need to implement any more code, so sit back and execute the code blocks! Some code comments are provided if you find yourself interested in the functionality.

In [54]:

```
def learning_curves(X_train, y_train, X_test, y_test):
    """ Calculates the performance of several models with varying sizes of
    training data.
    The learning and testing error rates for each model are then plotted. """

    print "Creating learning curve graphs for max_depths of 1, 3, 6, and 10. . ."

    # Create the figure window
    fig = plt.figure(figsize=(10,8))

    # We will vary the training set size so that we have 50 different sizes
    sizes = np.rint(np.linspace(1, len(X_train), 50)).astype(int)
    train_err = np.zeros(len(sizes))
    test_err = np.zeros(len(sizes))

    # Create four different models based on max_depth
    for k, depth in enumerate([1,3,6,10]):

        for i, s in enumerate(sizes):

            # Setup a decision tree regressor so that it learns a tree with
            max_depth = depth
            regressor = DecisionTreeRegressor(max_depth = depth)

            # Fit the learner to the training data
            regressor.fit(X_train[:s], y_train[:s])

            # Find the performance on the training set
            train_err[i] = performance_metric(y_train[:s], regressor.predict(X_train[:s]))

            # Find the performance on the testing set
            test_err[i] = performance_metric(y_test, regressor.predict(X_test))

        # Subplot the learning curve graph
        ax = fig.add_subplot(2, 2, k+1)
        ax.plot(sizes, test_err, lw = 2, label = 'Testing Error')
        ax.plot(sizes, train_err, lw = 2, label = 'Training Error')
        ax.legend()
        ax.set_title('max_depth = %s'%(depth))
        ax.set_xlabel('Number of Data Points in Training Set')
        ax.set_ylabel('Total Error')
        ax.set_xlim([0, len(X_train)])

    # Visual aesthetics
    fig.suptitle('Decision Tree Regressor Learning Performances', fontsize=18, y=1.03)
    fig.tight_layout()
    fig.show()
```

In [55]:

```
def model_complexity(X_train, y_train, X_test, y_test):  
    """ Calculates the performance of the model as model complexity increases.  
        The learning and testing errors rates are then plotted. """  
  
    print "Creating a model complexity graph. . . "  
  
    # We will vary the max_depth of a decision tree model from 1 to 14  
    max_depth = np.arange(1, 14)  
    train_err = np.zeros(len(max_depth))  
    test_err = np.zeros(len(max_depth))  
  
    for i, d in enumerate(max_depth):  
        # Setup a Decision Tree Regressor so that it learns a tree with depth d  
        regressor = DecisionTreeRegressor(max_depth = d)  
  
        # Fit the learner to the training data  
        regressor.fit(X_train, y_train)  
  
        # Find the performance on the training set  
        train_err[i] = performance_metric(y_train, regressor.predict(X_train))  
  
        # Find the performance on the testing set  
        test_err[i] = performance_metric(y_test, regressor.predict(X_test))  
  
    # Plot the model complexity graph  
    pl.figure(figsize=(7, 5))  
    pl.title('Decision Tree Regressor Complexity Performance')  
    pl.plot(max_depth, test_err, lw=2, label = 'Testing Error')  
    pl.plot(max_depth, train_err, lw=2, label = 'Training Error')  
    pl.legend()  
    pl.xlabel('Maximum Depth')  
    pl.ylabel('Total Error')  
    pl.show()
```

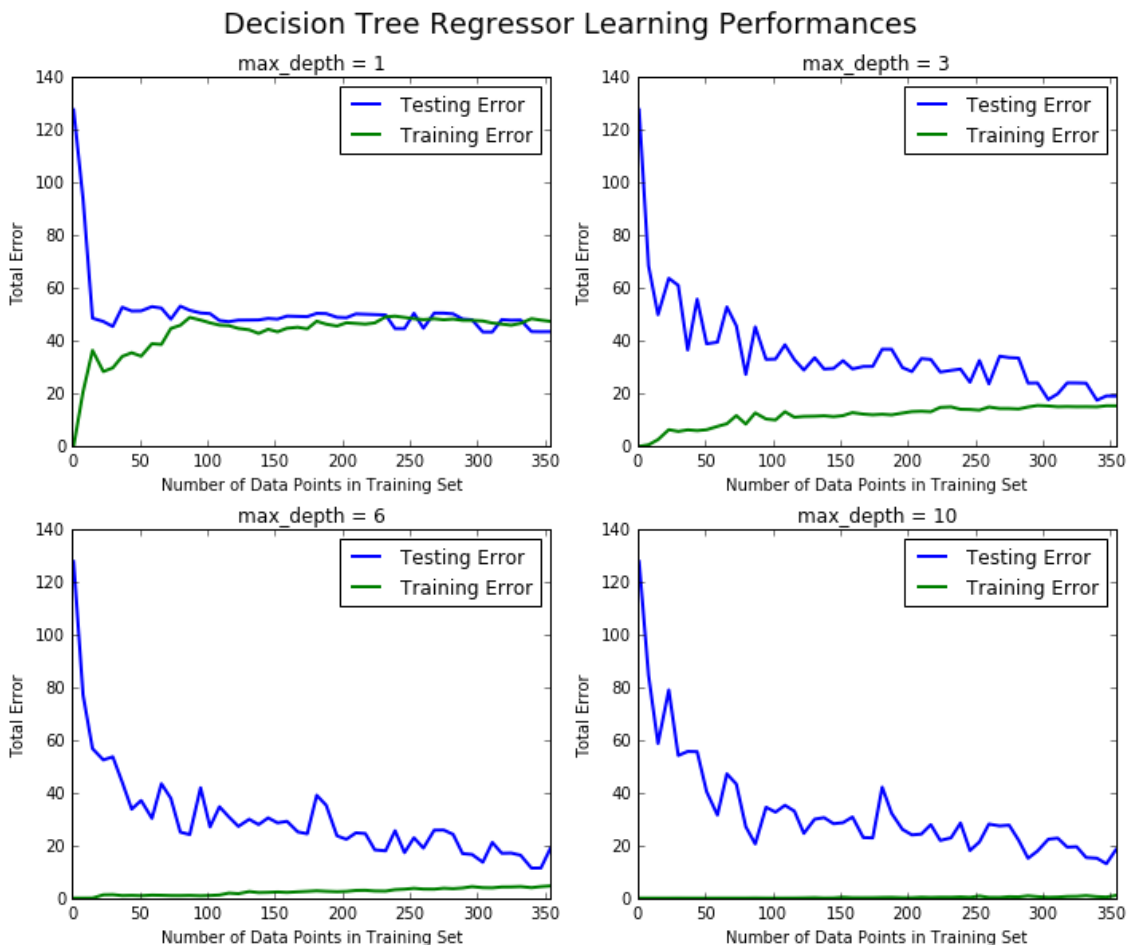
Analyzing Model Performance

In this third section of the project, you'll take a look at several models' learning and testing error rates on various subsets of training data. Additionally, you'll investigate one particular algorithm with an increasing `max_depth` parameter on the full training set to observe how model complexity affects learning and testing errors. Graphing your model's performance based on varying criteria can be beneficial in the analysis process, such as visualizing behavior that may not have been apparent from the results alone.

In [56]:

```
learning_curves(X_train, y_train, X_test, y_test)
```

Creating learning curve graphs for max_depths of 1, 3, 6, and 10. . .



Question 7

Choose one of the learning curve graphs that are created above. What is the max depth for the chosen model? As the size of the training set increases, what happens to the training error? What happens to the testing error?

I choose Max-depth of 6. As the size of training set increases, testing error are decreased slowly because of there are more data for training and resolved high biased of data. Training error will be slowly increase because of the model won't be able to fit all the training data that increase over time. We can describe the data in two phase. First phase is the underfited data, when the data size is not enough for training, the error rate of testing will be high because the system don't have enough data to learn. In the meanwhile, the more data is in the training set, the system can train and get more precision.

Question 8

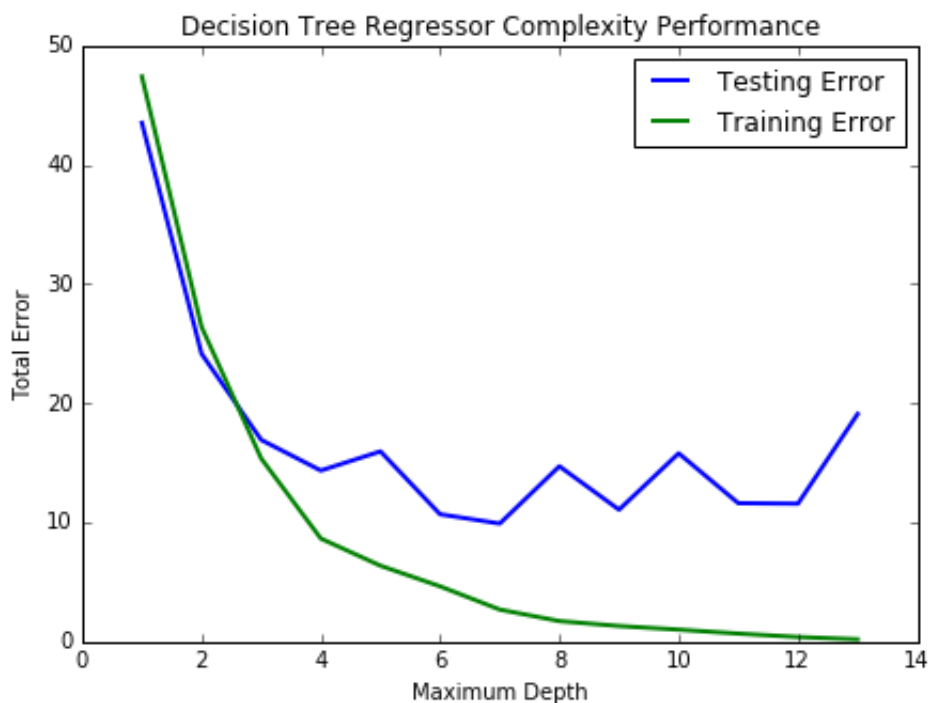
Look at the learning curve graphs for the model with a max depth of 1 and a max depth of 10. When the model is using the full training set, does it suffer from high bias or high variance when the max depth is 1? What about when the max depth is 10?

max depth of 1 are high bias in the first phase training error will increase dramatically because of the underfit of the data for training. The best performance is to get training error and testing error converge as close as possible while slowly converge which is max depth of 1 is not. while max depth of 10 suffer from high variance where even we increase the amount of data, training error and testing error won't converge anymore than this because of the overfit of the data that make traninig error are too overfit and no error at all for training.

In [57]:

```
model_complexity(X_train, y_train, X_test, y_test)
```

Creating a model complexity graph. . .



Question 9

From the model complexity graph above, describe the training and testing errors as the max depth increases. Based on your interpretation of the graph, which max depth results in a model that best generalizes the dataset? Why?

while max depth increase, the training error and test error will reduce but at one point where max depth increasing will not affected testing error anymore because of high variance of the data which training error is reducing while testing error won't follow along. So max-depth of 6 will be the best generalizes result for this dataset because of max-depth of 6 where testing error won't be reduced anymore

Model Prediction

In this final section of the project, you will make a prediction on the client's feature set using an optimized model from `fit_model`. When applying grid search along with cross-validation to optimize your model, it would typically be performed and validated on a training set and subsequently evaluated on a **dedicated test set**. In this project, the optimization below is performed on the *entire dataset* (as opposed to the training set you made above) due to the many outliers in the data. Using the entire dataset for training provides for a less volatile prediction at the expense of not testing your model's performance.

To answer the following questions, it is recommended that you run the code blocks several times and use the median or mean value of the results.

Question 10

Using grid search on the entire dataset, what is the optimal `max_depth` parameter for your model? How does this result compare to your intial intuition?

Hint: Run the code block below to see the max depth produced by your optimized model.

In [58]:

```
print "Final model has an optimal max_depth parameter of", reg.get_params()  
[ 'max_depth' ]
```

Final model has an optimal max_depth parameter of 6

the Optimal max-depth is 6 which is correct as same as the graph describe and my intuition. After max-depth of 6, testing error are not declined anymore but increaseing while max-depth is greater. the max-depth of 4 is nearly one of my choice also but unfortunately when I was looked in to max-depth of 6 where testing error is a little bit lower but training error is dramatic lower. So I choose max-depth of 6 instead

Question 11

With your parameter-tuned model, what is the best selling price for your client's home? How does this selling price compare to the basic statistics you calculated on the dataset?

Hint: Run the code block below to have your parameter-tuned model make a prediction on the client's home.

In [59]:

```
sale_price = reg.predict(CLIENT_FEATURES)
print "Predicted value of client's home: {:.3f}".format(sale_price[0])
```

Predicted value of client's home: 20.766

Very precise. from median house price which is 21.2, error rate just around 2%

Question 12 (Final Question):

In a few sentences, discuss whether you would use this model or not to predict the selling price of future clients' homes in the Greater Boston area.

Not sure until we try other predict model for this dataset. but in my opinion error rate for 2% is very acceptable