

Machine Learning Engineer Nanodegree

Model Evaluation & Validation

Project 1: Predicting Boston Housing Prices

Welcome to the first project of the Machine Learning Engineer Nanodegree! In this notebook, some template code has already been written. You will need to implement additional functionality to successfully answer all of the questions for this project. Unless it is requested, do not modify any of the code that has already been included. In this template code, there are four sections which you must complete to successfully produce a prediction with your model. Each section where you will write code is preceded by a **STEP X** header with comments describing what must be done. Please read the instructions carefully!

In addition to implementing code, there will be questions that you must answer that relate to the project and your implementation. Each section where you will answer a question is preceded by a **QUESTION X** header. Be sure that you have carefully read each question and provide thorough answers in the text boxes that begin with **"Answer:"**. Your project submission will be evaluated based on your answers to each of the questions.

A description of the dataset can be found [here \(https://archive.ics.uci.edu/ml/datasets/Housing\)](https://archive.ics.uci.edu/ml/datasets/Housing), which is provided by the **UCI Machine Learning Repository**.

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 [23]:

```
# 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
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 [24]:

```
# Number of houses in the dataset
total_houses = housing_features.shape[0]

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

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

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

# 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!

Answer: 1) CRIM: per capita crime rate by town - This is the number of crimes reported per person in town. Calculated by dividing the total number of reported criminal cases by the total number of persons in the town.

2) NOX: nitric oxides concentration(parts per 100 million) - This is a measure of parts of nitric oxide per 100 million parts of air.

3) DIS: weighted distances to five boston employment centres - The weighted distance is measured by taking the square root of the summation of the squares of the distances from each employment centre (each squared distance being multiplied by its weight as per it's significance).The weights are used to assign greater importance to the distance of one employment centre over the other.

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 [25]:

```
print CLIENT_FEATURES
```

```
[[11.95, 0.0, 18.1, 0, 0.659, 5.609, 90.0, 1.385, 24, 680.0, 20.2, 33  
2.09, 12.13]]
```

Answer: CRIM = 11.95 NOX = 0.659 DIS = 1.385

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 [26]:

```
# 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, X_test, y_train, y_test = train_test_split(X,y,test_size=0.3,random_s
tate=0)

    # 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, housin
g_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?

Answer: The data that we have access to needs to be split into training and testing data so that we can estimate the performance of our model against the test data after the training is done. It will help in preventing over-fitting to the training data. This cross validation will help us improve and generalize our model to give better results on unknown data.

We do need to make sure that the training and testing data have a good mix of different data points so as to have a robust model.

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 [27]:

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

def performance_metric(y_true, y_predict):
    """ Calculates and returns the total error between true and predicted values
        based on a performance metric chosen by the student. """

    error = mean_squared_error(y_true, y_predict)
    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)

Answer: Predicting housing prices is a case of regression. So, we have two choices 1) Mean Squared Error (MSE) 2) Mean Absolute Error (MAE)

The choice between these two depends on the application. I have chosen Mean Squared Error as the metric for Boston Housing because MSE would penalize the model more in case of large errors. As we want our price prediction to be close to actual prices it would be better that we emphasize large errors so that we can tune the model better. MSE would emphasize larger errors better than MAE. MSE is differentiable and we could find its minima by differentiating it and equating to zero, thus giving us the best parameter value for our model.

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 2**. 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 [28]:

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

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()

    # 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(mean_squared_error, greater_is_better=False)

    # Make the GridSearchCV object
    reg = GridSearchCV(regressor, parameters, scoring_function, cv=10)

    # Fit the Learner to the data to obtain the optimal model with tuned parameter
    reg.fit(X, y)

    # 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?

Answer: Grid search algorithm is used in cross validation to tune the parameter values, that are not directly learnt within estimators, to find the best combination for the estimator. The algorithm goes through multiple combinations of parameter tunes, cross-validating as it goes to determine which tune gives the best performance.

Grid Search consists of: 1) an estimator (regressor or classifier) 2) a parameter space 3) a method for searching or sampling candidates 4) a cross-validation scheme 5) a score function

Sklearn's GridSearchCV does an exhaustive search over specified parameter space for the estimator. It generates all the possible parameter combinations and then evaluates its score by using cross validation to decide the best combination. An example application would be tuning the parameters of an SVM while going through cross validation iterations to get the best fit for the 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?

Answer: Cross Validation is a model evaluation method. The data partitioned into training and testing sets is used to train and then validate the model. We have access to limited amount of data to build our model that needs to be partitioned into training and testing sets. We want to maximize our training as well testing sets so as to build a robust model.

Cross validation has various iterators like the K Fold where the data is divided into K folds. K-1 fold is used for training the model and the remaining 1 fold is used to validate it. This iteration goes for K times and performance metric is averaged over it.

A grid search algorithm needs a performance metric which is measured by cross validation on the training/validation set. Cross validation will help in identifying the parameter combination picked from the parameter space that gives the best performing and robust model.

In step 4, I have used 10 fold cross validation for GridSearchCV. It has been noted in a research paper [1] by Kohavi, Ron that stratified 10 fold cross validation would be the most appropriate choice for real world applications. The right choice for K might vary from one requirement to other but as a thumb rule I have taken 10 as my choice for step 4.

As K increases the bias reduces but the variance increases. Also, a larger K is more expensive. As our training/testing dataset is small, in order to reduce bias and also keep variance in check, 10 fold cross validation seems like an appropriate choice. With the default 3 fold CV, the data would have been divided into 1/3 for testing and 2/3 for training for the 3 iterations and with a relatively small sample like Boston Housing, I think this will introduce bias.

[1] Kohavi, Ron. A study of cross-validation and bootstrap for accuracy estimation and model selection. In IJCAI, 1995.

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 [29]:

```
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 [30]:

```
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

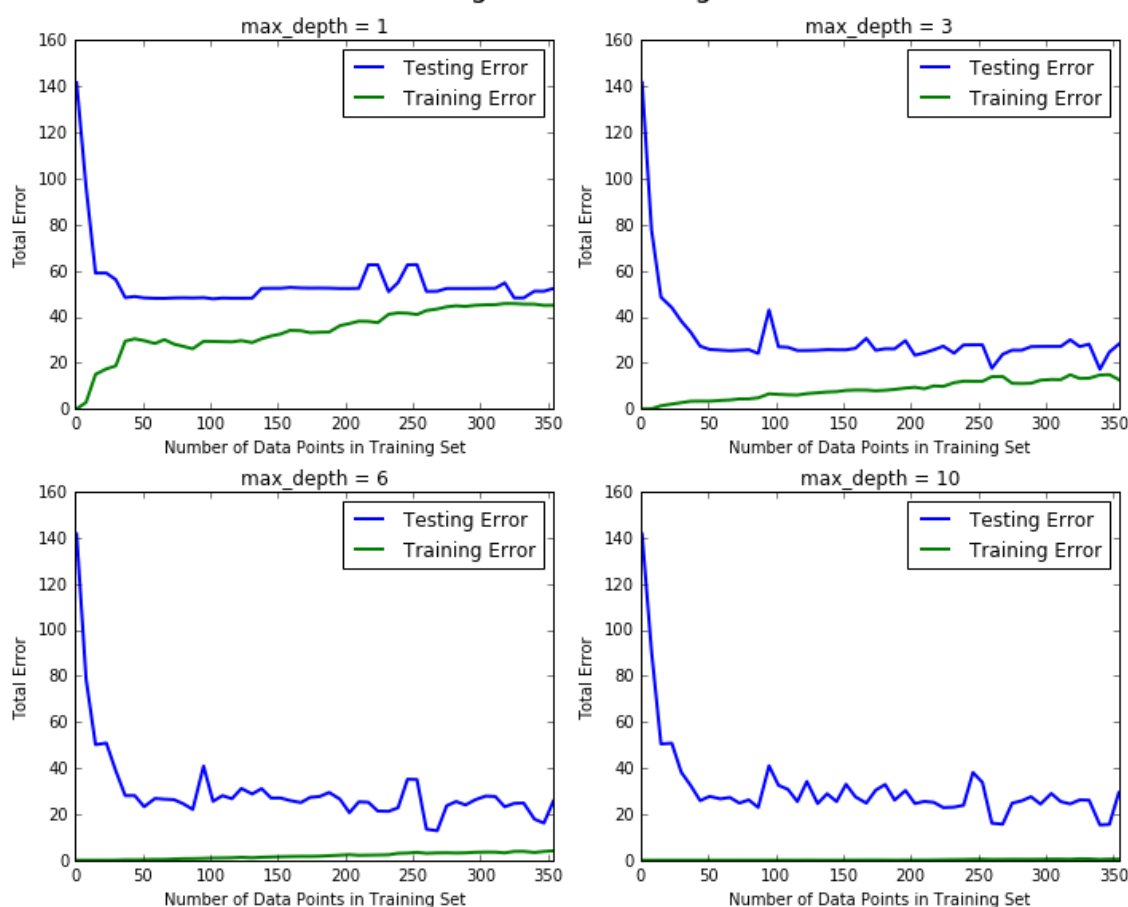
In this third section of the project, you'll take a look at several models' learning and testing error rates on various subsets of training data. Additionally, you'll investigate one particular algorithm with an increasing `max_depth` parameter on the full training set to observe how model complexity 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 [31]:

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

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

Decision Tree Regressor Learning Performances



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?

Answer: Considering the 2nd learning curve for the model having `max_depth = 3`.

The training error increases almost linearly till the number of training points reach 260, after that training error seems to be becoming constant. This indicates that the model has reached the point where increasing the data points will have no significant impact on its training performance. Also, the training error settles at a significantly high error suggesting that the model has high bias.

The testing error drops exponentially as we increase the data points from 0 to 50 indicating that the model has shown significant improvement as we have increased the data points. After this the testing error becomes nearly constant and settles at a significantly high error. This means that the model performance does not improve much even if we keep on increasing the data points beyond 50.

Both the training and testing error converge to nearly the same error and the model shows high bias.

As the size of the training set increases, there is an increase in the training error suggesting that it becomes increasingly difficult for the model to fit to the training data as its size increases. The testing error decreases with increase in training set suggesting that the model becomes more efficient in generalizing the test data.

As the model complexity increases, the training error decreases for a given number of training data points (say number of data points = 200) suggesting we get a better fit for our training set. The testing error however is seen to become stagnant as the `model_depth` reaches 6. This means that at `model_depth` greater than 6 we get a model that overfits to the training data.

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?

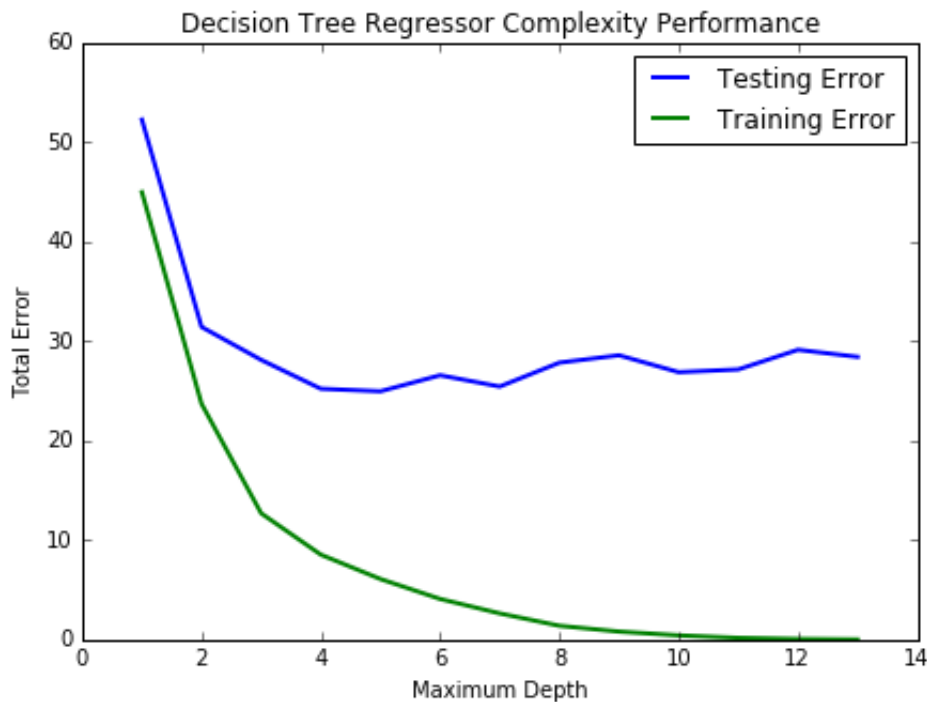
Answer: At `max_depth 1`, both the training and testing error converge at a high value as we use the full training set. As the training error is high, the model suffers from high bias i.e. the model is not able to represent underlying relationship of the data.

At `max_depth = 10`, the model suffers from high variance as the training error is almost negligible whereas we have a greater testing error. This model overfits the data to training set.

In [32]:

```
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?

Answer: The training error falls rapidly as the max_depth (model complexity) increases and becomes negligible as max_depth reaches 10. The testing error decreases till max_depth 5 and then starts to increase slightly.

The graph shows that the model starts overfitting the data at max_depth = 6 and above because the training error becomes very small but the testing error does not show any improvement, there is high variance.

The max_depth of 4 or 5 seems to best generalize the dataset as the testing error reaches a minimum at these levels. After these levels the increase in the model complexity does not help in generalizing the data as the model starts overfitting to test data.

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 initial intuition?

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

In [33]:

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

Final model has an optimal max_depth parameter of 5

Answer: The optimal max_depth of 5 for this model does match with my initial intuition.

I executed the code blocks from step 4 onwards several times and got varying optimal max_depth parameters from 4 to 10 with each run. The median max_depth was 6 and the mean was 6.3, however 5 seems to be a better choice.

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 [34]:

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

Predicted value of client's home: 20.968

Answer: The best selling price for the client's home is 20.968 based on the optimized model. The selling price is near to the median price of 21.2 and not very far from the mean price of 22.533 for the boston dataset.

The mean price for several runs of the code was 20.46 and the median was 20.776.

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.

Answer: The best optimized model also has a large mean squared error on test sets and so the prediction is likely to be incorrect. I would not use this model to predict selling price of future client's home. There was no improvement in testing errors as the number of data points were increased beyond 50. So, adding more data points would also not help. The decision tree regressor is quite robust to outliers, so the outliers would not have impacted the model performance.

Random forest regressor or SVMs could have been other choices that we could have used for better results.