

boston_housing

January 5, 2018

1 Machine Learning Engineer Nanodegree

1.1 Model Evaluation & Validation

1.2 Project: Predicting Boston Housing Prices

Welcome to the first project of the Machine Learning Engineer Nanodegree! In this notebook, some template code has already been provided for you, and you will need to implement additional functionality to successfully complete this project. You will not need to modify the included code beyond what is requested. Sections that begin with '**Implementation**' in the header indicate that the following block of code will require additional functionality which you must provide. Instructions will be provided for each section and the specifics of the implementation are marked in the code block with a 'TODO' statement. Please be sure to read the instructions carefully!

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

Note: Code and Markdown cells can be executed using the **Shift + Enter** keyboard shortcut. In addition, Markdown cells can be edited by typically double-clicking the cell to enter edit mode.

1.3 Getting Started

In this project, you will evaluate the performance and predictive power of a model that has been trained and tested on data collected from homes in suburbs of Boston, Massachusetts. A model trained on this data that is seen as a *good fit* could then be used to make certain predictions about a home — in particular, its monetary value. This model would prove to be invaluable for someone like a real estate agent who could make use of such information on a daily basis.

The dataset for this project originates from the [UCI Machine Learning Repository](#). The Boston housing data was collected in 1978 and each of the 506 entries represent aggregated data about 14 features for homes from various suburbs in Boston, Massachusetts. For the purposes of this project, the following preprocessing steps have been made to the dataset: - 16 data points have an 'MEDV' value of 50.0. These data points likely contain **missing or censored values** and have been removed. - 1 data point has an 'RM' value of 8.78. This data point can be considered an

outlier and has been removed. - The features 'RM', 'LSTAT', 'PTRATIO', and 'MEDV' are essential. The remaining **non-relevant features** have been excluded. - The feature 'MEDV' has been **multiplicatively scaled** to account for 35 years of market inflation.

Run the code cell below to load the Boston housing dataset, along with a few of the necessary Python libraries required for this project. You will know the dataset loaded successfully if the size of the dataset is reported.

```
In [1]: # Import libraries necessary for this project
import numpy as np
import pandas as pd
from sklearn.cross_validation import ShuffleSplit

# Import supplementary visualizations code visuals.py
import visuals as vs

# Pretty display for notebooks
%matplotlib inline

# Load the Boston housing dataset
data = pd.read_csv('housing.csv')
prices = data['MEDV']
features = data.drop('MEDV', axis = 1)

# Success
print "Boston housing dataset has {} data points with {} variables each.".format(*data
```

C:\Users\Tron\Anaconda2\lib\site-packages\sklearn\cross_validation.py:44: DeprecationWarning: This module will be removed in 0.20.", DeprecationWarning)

Boston housing dataset has 489 data points with 4 variables each.

C:\Users\Tron\Anaconda2\lib\site-packages\sklearn\learning_curve.py:23: DeprecationWarning: This module will be removed in 0.20.", DeprecationWarning)

1.4 Data Exploration

In this first section of this project, you will make a cursory investigation about the Boston housing data and provide your observations. Familiarizing yourself with the data through an explorative process is a fundamental practice to help you better understand and justify your results.

Since the main goal of this project is to construct a working model which has the capability of predicting the value of houses, we will need to separate the dataset into **features** and the **target variable**. The **features**, 'RM', 'LSTAT', and 'PTRATIO', give us quantitative information about each data point. The **target variable**, 'MEDV', will be the variable we seek to predict. These are stored in features and prices, respectively.

1.4.1 Implementation: Calculate Statistics

For your very first coding implementation, you will calculate descriptive statistics about the Boston housing prices. Since numpy has already been imported for you, use this library to perform the necessary calculations. These statistics will be extremely important later on to analyze various prediction results from the constructed model.

In the code cell below, you will need to implement the following: - Calculate the minimum, maximum, mean, median, and standard deviation of 'MEDV', which is stored in prices. - Store each calculation in their respective variable.

```
In [2]: # TODO: Minimum price of the data
        minimum_price = np.amin(prices)

        # TODO: Maximum price of the data
        maximum_price = np.amax(prices)

        # TODO: Mean price of the data
        mean_price = np.mean(prices)

        # TODO: Median price of the data
        median_price = np.median(prices)

        # TODO: Standard deviation of prices of the data
        std_price = np.std(prices)

        import scipy
        skewness_price = scipy.stats.skew(prices)

        # Show the calculated statistics
        print "Statistics for Boston housing dataset:\n"
        print "Minimum price: ${:,.2f}".format(minimum_price)
        print "Maximum price: ${:,.2f}".format(maximum_price)
        print "Mean price: ${:,.2f}".format(mean_price)
        print "Median price ${:,.2f}".format(median_price)
        print "Standard deviation of prices: ${:,.2f}".format(std_price)
        print "Skewness of prices: {:.2f}".format(skewness_price)
```

Statistics for Boston housing dataset:

```
Minimum price: $105,000.00
Maximum price: $1,024,800.00
Mean price: $454,342.94
Median price $438,900.00
Standard deviation of prices: $165,171.13
Skewness of prices: 0.77
```

1.4.2 Question 1 - Feature Observation

As a reminder, we are using three features from the Boston housing dataset: 'RM', 'LSTAT', and 'PTRATIO'. For each data point (neighborhood): - 'RM' is the average number of rooms among homes in the neighborhood. - 'LSTAT' is the percentage of homeowners in the neighborhood considered "lower class" (working poor). - 'PTRATIO' is the ratio of students to teachers in primary and secondary schools in the neighborhood.

*Using your intuition, for each of the three features above, do you think that an increase in the value of that feature would lead to an **increase** in the value of 'MEDV' or a **decrease** in the value of 'MEDV'? Justify your answer for each.*

Hint: Would you expect a home that has an 'RM' value of 6 be worth more or less than a home that has an 'RM' value of 7?

I think that RM is positively correlated with MEDV whereas LSTAT and PTRATIO are negatively correlated with MEDV. INTUITIONS : - For RM the idea is that globally a house with more rooms is likely to have a higher price - For LSTAT the idea is that lower class workers are unlikely to be owners of highly priced houses - For PTRATIO I think that area with rich inhabitants living in highly priced houses are unlikely to have schools with overpopulated classes

1.5 Developing a Model

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

1.5.1 Implementation: Define a Performance Metric

It is difficult to measure the quality of a given model without quantifying its performance over training and testing. This is typically done using some type of performance metric, whether it is through calculating some type of error, the goodness of fit, or some other useful measurement. For this project, you will be calculating the *coefficient of determination*, R2, to quantify your model's performance. The coefficient of determination for a model is a useful statistic in regression analysis, as it often describes how "good" that model is at making predictions.

The values for R2 range from 0 to 1, which captures the percentage of squared correlation between the predicted and actual values of the **target variable**. A model with an R2 of 0 is no better than a model that always predicts the *mean* of the target variable, whereas a model with an R2 of 1 perfectly predicts the target variable. Any value between 0 and 1 indicates what percentage of the target variable, using this model, can be explained by the **features**. *A model can be given a negative R2 as well, which indicates that the model is **arbitrarily worse** than one that always predicts the mean of the target variable.*

For the `performance_metric` function in the code cell below, you will need to implement the following: - Use `r2_score` from `sklearn.metrics` to perform a performance calculation between `y_true` and `y_predict`. - Assign the performance score to the `score` variable.

```
In [3]: # TODO: Import 'r2_score'
        from sklearn.metrics import r2_score
```

```
def performance_metric(y_true, y_predict):
    """ Calculates and returns the performance score between
        true and predicted values based on the metric chosen. """

    # TODO: Calculate the performance score between 'y_true' and 'y_predict'
    score = r2_score(y_true, y_predict)

    # Return the score
    return score
```

1.5.2 Question 2 - Goodness of Fit

Assume that a dataset contains five data points and a model made the following predictions for the target variable:

True Value	Prediction
3.0	2.5
-0.5	0.0
2.0	2.1
7.0	7.8
4.2	5.3

Would you consider this model to have successfully captured the variation of the target variable? Why or why not?

Run the code cell below to use the `performance_metric` function and calculate this model's coefficient of determination.

```
In [4]: # Calculate the performance of this model
score = performance_metric([3, -0.5, 2, 7, 4.2], [2.5, 0.0, 2.1, 7.8, 5.3])
print "Model has a coefficient of determination, R^2, of {:.3f}.".format(score)
```

Model has a coefficient of determination, R², of 0.923.

I think the model has successfully captured the variation of the dataset because the R²_score is close to 1 and it is a good performance with only 5 points in a dataset.

1.5.3 Implementation: Shuffle and Split Data

Your next implementation requires that you take the Boston housing dataset and split the data into training and testing subsets. Typically, the data is also shuffled into a random order when creating the training and testing subsets to remove any bias in the ordering of the dataset.

For the code cell below, you will need to implement the following: - Use `train_test_split` from `sklearn.cross_validation` to shuffle and split the features and prices data into training and testing sets. - Split the data into 80% training and 20% testing. - Set the `random_state` for `train_test_split` to a value of your choice. This ensures results are consistent. - Assign the train and testing splits to `X_train`, `X_test`, `y_train`, and `y_test`.

```
In [5]: # TODO: Import 'train_test_split'
        from sklearn.cross_validation import train_test_split

        # TODO: Shuffle and split the data into training and testing subsets
        X_train, X_test, y_train, y_test = train_test_split(features, prices, test_size=0.2, r

        # Success
        print "Training and testing split was successful."
```

Training and testing split was successful.

1.5.4 Question 3 - Training and Testing

What is the benefit to splitting a dataset into some ratio of training and testing subsets for a learning algorithm?

Hint: What could go wrong with not having a way to test your model?

It gives some kind of feedback on the robustness of the algorithm. It allows to measure whether the model generalizes well or not on new examples and to readjust the model if necessary.

1.6 Analyzing Model Performance

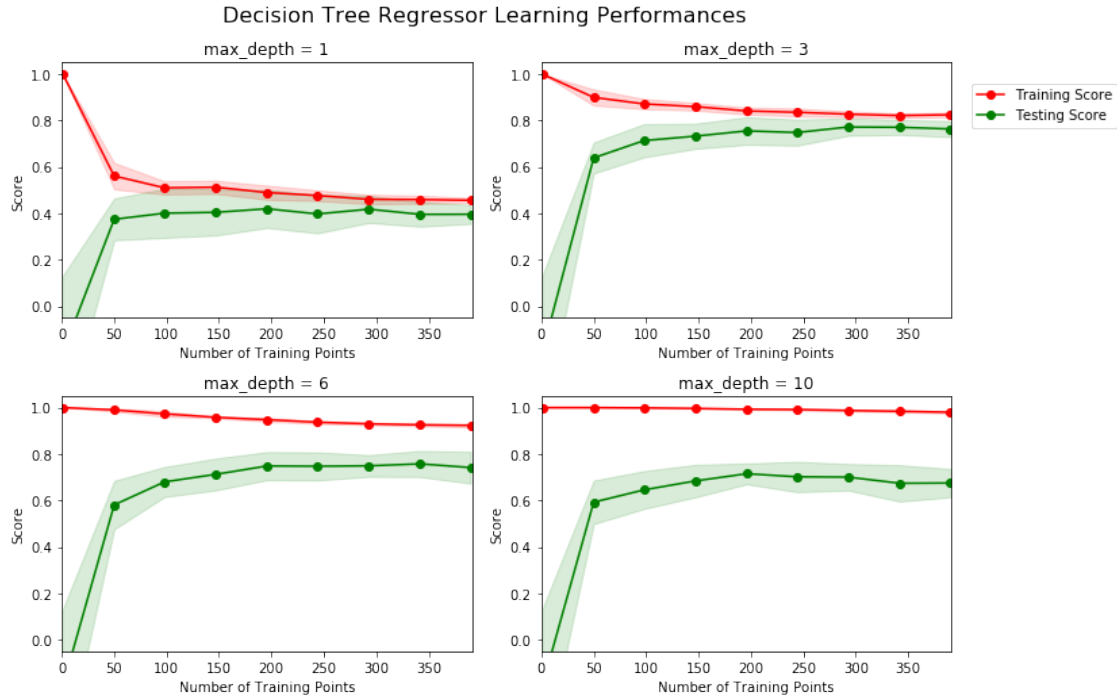
In this third section of the project, you'll take a look at several models' learning and testing performances on various subsets of training data. Additionally, you'll investigate one particular algorithm with an increasing 'max_depth' parameter on the full training set to observe how model complexity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

1.6.1 Learning Curves

The following code cell produces four graphs for a decision tree model with different maximum depths. Each graph visualizes the learning curves of the model for both training and testing as the size of the training set is increased. Note that the shaded region of a learning curve denotes the uncertainty of that curve (measured as the standard deviation). The model is scored on both the training and testing sets using R2, the coefficient of determination.

Run the code cell below and use these graphs to answer the following question.

```
In [6]: # Produce learning curves for varying training set sizes and maximum depths
        vs.ModelLearning(features, prices)
```



1.6.2 Question 4 - Learning the Data

Choose one of the graphs above and state the maximum depth for the model. What happens to the score of the training curve as more training points are added? What about the testing curve? Would having more training points benefit the model?

Hint: Are the learning curves converging to particular scores?

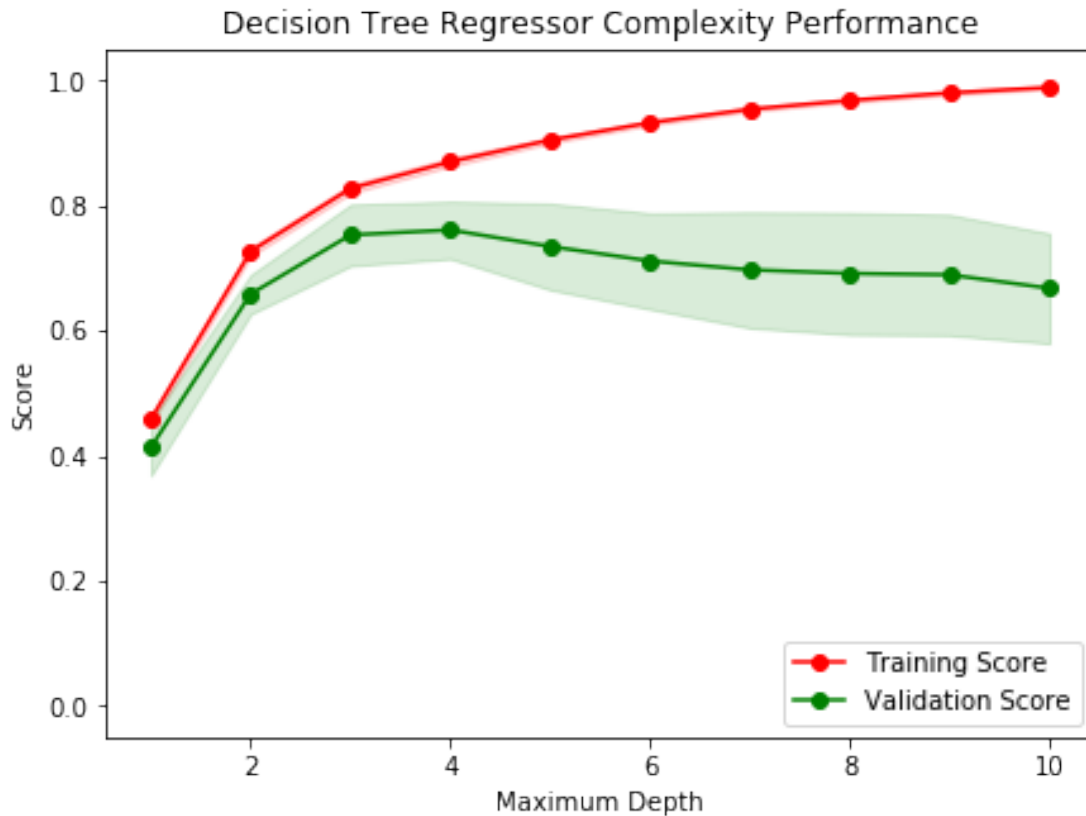
Top-right graph `max_depth = 3` As more training points are added, the training score is decreasing and seems to converge to a value around 0.825. As more training points are added, the testing score is increasing and seems to converge to a value around 0.75. More training points are absolutely necessary to more complex models such as `max_depth = 6` and `max_depth = 10` models. I am not sure the benefit of having more examples would be significant and would be worth the price to obtain them, for a chosen low-complexity model such as the `max_depth = 3` one.

1.6.3 Complexity Curves

The following code cell produces a graph for a decision tree model that has been trained and validated on the training data using different maximum depths. The graph produces two complexity curves — one for training and one for validation. Similar to the **learning curves**, the shaded regions of both the complexity curves denote the uncertainty in those curves, and the model is scored on both the training and validation sets using the `performance_metric` function.

Run the code cell below and use this graph to answer the following two questions.

```
In [7]: vs.ModelComplexity(X_train, y_train)
```



1.6.4 Question 5 - Bias-Variance Tradeoff

When the model is trained with a maximum depth of 1, does the model suffer from high bias or from high variance? How about when the model is trained with a maximum depth of 10? What visual cues in the graph justify your conclusions?

Hint: How do you know when a model is suffering from high bias or high variance?

max_depth = 1 high bias As seen on the complexity performance graph, the training score and the testing score are both low (validation score around 0.4, training score around 0.45), which is due to underfitting i.e. the fact that the model is not complex enough to learn well the data so it has a low score both on new examples (testing set) and on examples it has already 'seen' (training set).

max_depth = 10 high variance As seen on the complexity performance graph, the algorithm performs well on the training set (training score close to 1) but do not generalizes well on new examples (testing score close to 0.6). This is a case of overfitting i.e. the algorithm learns random noise from the data due to a too complex model.

1.6.5 Question 6 - Best-Guess Optimal Model

Which maximum depth do you think results in a model that best generalizes to unseen data? What intuition lead you to this answer?

The `max_depth = 4` model is the model that best generalizes to unseen data. This model has the best testing score on the complexity performance graph.

1.7 Evaluating Model Performance

In this final section of the project, you will construct a model and make a prediction on the client's feature set using an optimized model from `fit_model`.

1.7.1 Question 7 - Grid Search

What is the grid search technique and how it can be applied to optimize a learning algorithm?

It is a technique to optimize at the same time, with a few lines of code, many models and many parameters.

The hyperparameters and the types of models are chosen by the programmer and generally specified in a dictionary.

This technique will train all the combinations of hyperparameters and model types and compute a performance score according to an evaluation metrics (R^2 in this case). Then the algorithm will chose the combination that performs best.

1.7.2 Question 8 - Cross-Validation

What is the k-fold cross-validation training technique? What benefit does this technique provide for grid search when optimizing a model?

Hint: Much like the reasoning behind having a testing set, what could go wrong with using grid search without a cross-validated set?

It is a technique that consist of successively splitting the database and training the algorithm with different repartitions. For instance repeating k times the following steps: - split the database in k subset and chose the k -th subset as the testing set and the other subset as the training set - optimize the algorithm

The score is the average of the k scores obtained through the k iterations.

The benefit of this method is that all the data participate in the learning process, so the algorithm is likely to be more robust. So the ratio (quantity of data in the testing set / quantity of data in the entire dataset) is not as important as in the simpleholdout method.

As a consequence, the hyperparameters chosen by this method are more likely to be the best ones. The simple holdout method can lead to chose others parameters depending on how the splitis made. For instance, in my first submission of the project, I have mistakenly split the dataset with a `test_size` ratio of 0.25 "Implementation: Shuffle and Split Data" which has led the program to choose the `max_depth = 3` model instead of the `max_depth = 4` model in the Question 9 - Optimal Model.

1.7.3 Implementation: Fitting a Model

Your final implementation requires that you bring everything together and train a model using the **decision tree algorithm**. To ensure that you are producing an optimized model, you will train the model using the grid search technique to optimize the '`max_depth`' parameter for the decision tree. The '`max_depth`' parameter can be thought of as how many questions the decision

tree algorithm is allowed to ask about the data before making a prediction. Decision trees are part of a class of algorithms called *supervised learning algorithms*.

In addition, you will find your implementation is using `ShuffleSplit()` for an alternative form of cross-validation (see the `'cv_sets'` variable). While it is not the K-Fold cross-validation technique you describe in **Question 8**, this type of cross-validation technique is just as useful!. The `ShuffleSplit()` implementation below will create 10 (`'n_splits'`) shuffled sets, and for each shuffle, 20% (`'test_size'`) of the data will be used as the *validation set*. While you're working on your implementation, think about the contrasts and similarities it has to the K-fold cross-validation technique.

Please note that `ShuffleSplit` has different parameters in scikit-learn versions 0.17 and 0.18. For the `fit_model` function in the code cell below, you will need to implement the following: - Use `DecisionTreeRegressor` from `sklearn.tree` to create a decision tree regressor object. - Assign this object to the `'regressor'` variable. - Create a dictionary for `'max_depth'` with the values from 1 to 10, and assign this to the `'params'` variable. - Use `make_scorer` from `sklearn.metrics` to create a scoring function object. - Pass the `performance_metric` function as a parameter to the object. - Assign this scoring function to the `'scoring_fnc'` variable. - Use `GridSearchCV` from `sklearn.grid_search` to create a grid search object. - Pass the variables `'regressor'`, `'params'`, `'scoring_fnc'`, and `'cv_sets'` as parameters to the object. - Assign the `GridSearchCV` object to the `'grid'` variable.

```
In [8]: # TODO: Import 'make_scorer', 'DecisionTreeRegressor', and 'GridSearchCV'
        from sklearn.metrics import make_scorer
        from sklearn.tree import DecisionTreeRegressor
        from sklearn.grid_search import GridSearchCV
        #from sklearn import grid_search

        def fit_model(X, y):
            """ Performs grid search over the 'max_depth' parameter for a
                decision tree regressor trained on the input data [X, y]. """

            # Create cross-validation sets from the training data
            cv_sets = ShuffleSplit(X.shape[0], n_iter = 10, test_size = 0.20, random_state = 0)

            # TODO: Create a decision tree regressor object
            regressor = DecisionTreeRegressor()

            # TODO: Create a dictionary for the parameter 'max_depth' with a range from 1 to 10
            params = {'max_depth': range(1,11)}

            # TODO: Transform 'performance_metric' into a scoring function using 'make_scorer'
            scoring_fnc = make_scorer(performance_metric)

            # TODO: Create the grid search object
            grid = GridSearchCV(regressor, params, scoring = scoring_fnc, cv = cv_sets)

            # Fit the grid search object to the data to compute the optimal model
            grid = grid.fit(X, y)
```

```
# Return the optimal model after fitting the data
return grid.best_estimator_
```

C:\Users\Tron\Anaconda2\lib\site-packages\sklearn\grid_search.py:43: DeprecationWarning: This is a DeprecationWarning)

1.7.4 Making Predictions

Once a model has been trained on a given set of data, it can now be used to make predictions on new sets of input data. In the case of a *decision tree regressor*, the model has learned *what the best questions to ask about the input data are*, and can respond with a prediction for the **target variable**. You can use these predictions to gain information about data where the value of the target variable is unknown — such as data the model was not trained on.

1.7.5 Question 9 - Optimal Model

What maximum depth does the optimal model have? How does this result compare to your guess in *Question 6*?

Run the code block below to fit the decision tree regressor to the training data and produce an optimal model.

```
In [9]: # Fit the training data to the model using grid search
        reg = fit_model(X_train, y_train)

        # Produce the value for 'max_depth'
        print "Parameter 'max_depth' is {} for the optimal model.".format(reg.get_params()['max
```

Parameter 'max_depth' is 4 for the optimal model.

According to the program, "Parameter 'max_depth' is 4 for the optimal model". It confirms the guess in question 6.

1.7.6 Question 10 - Predicting Selling Prices

Imagine that you were a real estate agent in the Boston area looking to use this model to help price homes owned by your clients that they wish to sell. You have collected the following information from three of your clients:

Feature	Client 1	Client 2	Client 3
Total number of rooms in home	5 rooms	4 rooms	8 rooms
Neighborhood poverty level (as %)	17%	32%	3%
Student-teacher ratio of nearby schools	15-to-1	22-to-1	12-to-1

What price would you recommend each client sell his/her home at? Do these prices seem reasonable given the values for the respective features?

Hint: Use the statistics you calculated in the **Data Exploration** section to help justify your re-

sponse.

Run the code block below to have your optimized model make predictions for each client's home.

```
In [10]: # Produce a matrix for client data
        client_data = [[5, 17, 15], # Client 1
                        [4, 32, 22], # Client 2
                        [8, 3, 12]]  # Client 3

        # Show predictions
        for i, price in enumerate(reg.predict(client_data)):
            print "Predicted selling price for Client {}'s home: ${:,.2f}".format(i+1, price)

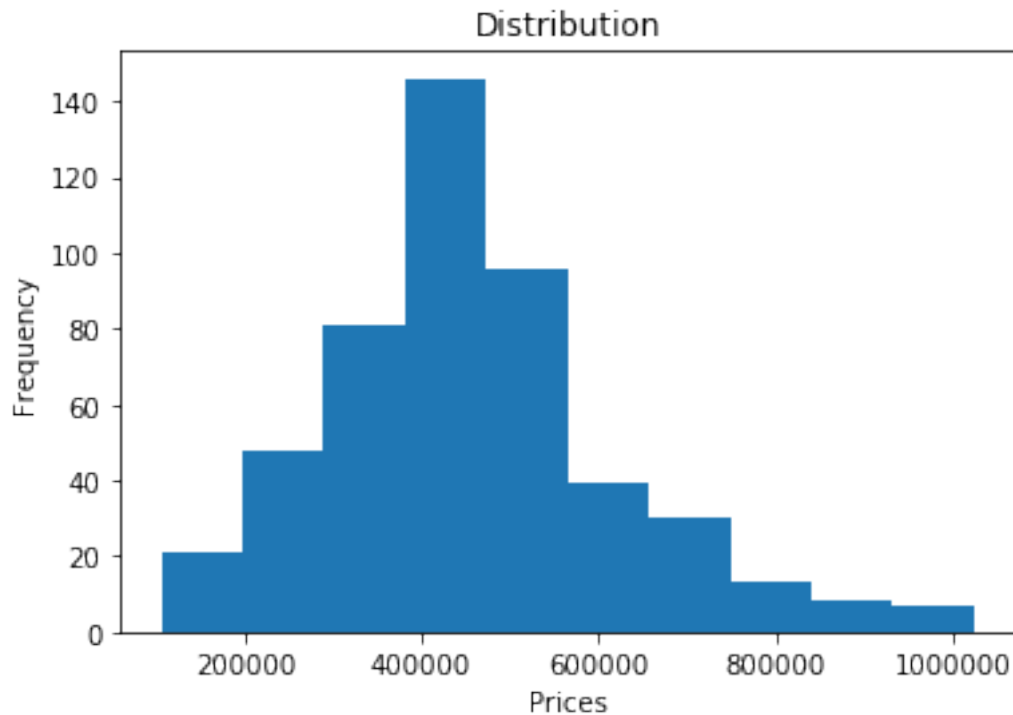
        import matplotlib.pyplot as plt
        plt.hist(prices)
        plt.title("Distribution")
        plt.xlabel("Prices")
        plt.ylabel("Frequency")
        plt.show()

        from IPython.display import Image, display
        from IPython.core.display import HTML
        a = Image(url= "https://udacity-github-sync-content.s3.amazonaws.com/_attachments/6500")
        b = Image(url= "https://udacity-github-sync-content.s3.amazonaws.com/_attachments/6500")
        display(a)
        display(b)
```

Predicted selling price for Client 1's home: \$403,025.00

Predicted selling price for Client 2's home: \$237,478.72

Predicted selling price for Client 3's home: \$931,636.36



<IPython.core.display.Image object>

<IPython.core.display.Image object>

I would suggest the prices given by the algorithm: Predicted selling price for Client 1's home: \$403,025.00 Predicted selling price for Client 2's home: \$237,478.72 Predicted selling price for Client 3's home: \$931,636.36

Those prices are coherent with the intuitions in question 1 (RM positively correlated with MEDV while LSTAT and PTRATIO negatively correlated with MEDV), according to the respective values of the features so the prices seems reasonable to me.

Statistics for Boston housing dataset:

Minimum price: \$105,000.00 Maximum price: \$1,024,800.00 Mean price: \$454,342.94 Median price \$438,900.00

House #3 graphs comments As shown on the histogram of prices, the house #3 has a predicted price (\$931,636.36) that is amongst the highest price range (the maximum price is \$1,024,800.00). It is coherent with point position in the 3 graphs: extreme right in the MEDV/RM graph (8 rooms), extreme left on both the MEDV/LSTAT and MEDV/PTRATIO graphs (3% of lower-class homeowners and a ratio of Students/Teacher of 12)

House #2 graphs comments As shown on the histogram of prices, the house #2 has a predicted price (\$237,478.72) that is amongst the second lowestest price range (Median price \$438,900.00). It is coherent with point position in the 3 graphs: extreme left in the MEDV/RM graph (4 rooms),

extreme right on both the MEDV/LSTAT and MEDV/PTRATIO graphs (32% of lower-class homeowners and a ratio of Students/Teacher of 22)

House #1 graphs comments As shown on the histogram of prices, the house #1 has a predicted price (\$403,025.00) that is amongst the middle price range (Minimum price: \$105,000.00). It is coherent with point position in the 3 graphs: close to the center/slightly on the left in the MEDV/RM graph (5 rooms), close to the center/slightly on the left on both the MEDV/LSTAT and MEDV/PTRATIO graphs (17% of lower-class home-owners and a ratio of Students/Teacher of 15)

1.7.7 Sensitivity

An optimal model is not necessarily a robust model. Sometimes, a model is either too complex or too simple to sufficiently generalize to new data. Sometimes, a model could use a learning algorithm that is not appropriate for the structure of the data given. Other times, the data itself could be too noisy or contain too few samples to allow a model to adequately capture the target variable — i.e., the model is underfitted. Run the code cell below to run the `fit_model` function ten times with different training and testing sets to see how the prediction for a specific client changes with the data it's trained on.

```
In [11]: vs.PredictTrials(features, prices, fit_model, client_data)
```

```
Trial 1: $391,183.33
Trial 2: $419,700.00
Trial 3: $415,800.00
Trial 4: $420,622.22
Trial 5: $418,377.27
Trial 6: $411,931.58
Trial 7: $399,663.16
Trial 8: $407,232.00
Trial 9: $351,577.61
Trial 10: $413,700.00
```

```
Range in prices: $69,044.61
```

1.7.8 Question 11 - Applicability

In a few sentences, discuss whether the constructed model should or should not be used in a real-world setting.

Hint: Some questions to answering: - *How relevant today is data that was collected from 1978?* - *Are the features present in the data sufficient to describe a home?* - *Is the model robust enough to make consistent predictions?* - *Would data collected in an urban city like Boston be applicable in a rural city?*

The model should not be used in the real world for several reasons: - It is not robust (Range in prices: \$73,357.39). - Moreover other features may impact the price of a house (beauty, material, energy efficiency, etc...). - Data collected in 1978 may not be relevant nowadays as Western nations experienced periods of price inflation the last decades (according to <http://www.usinflationcalculator.com/>, the cumulative rate of inflation between 1978 and 2017 is 273.9% in the U.S.) - The algorithm may not generalize very well if it is trained with only examples from Boston area

Note: Once you have completed all of the code implementations and successfully answered each question above, you may finalize your work by exporting the iPython Notebook as an HTML document. You can do this by using the menu above and navigating to

File -> Download as -> HTML (.html). Include the finished document along with this notebook as your submission.