# MLND - P1 - Vishakh Rayapeta

## 1) Statistical Analysis and Data Exploration

- Number of data points (houses)?

  *506*

- Number of features?

  *13*

- Minimum and maximum housing prices?

  *5.0, 50.0*

- Mean and median Boston housing prices?

  *22.533, 21.2*

- Standard deviation?

  *9.188*

## 2) Evaluating Model Performance

- Which measure of model performance is best to use for predicting Boston housing data and analyzing the errors? Why do you think this measurement most appropriate? Why might the other measurements not be appropriate here?

*The Decision Tree Regressor uses the 'Mean Squared Error' to measure the quality of the split. This is the reason to use this measure.*

- Why is it important to split the Boston housing data into training and testing data? What happens if you do not do this?

*The testing data is used to ensure that 'overfitting error' is not too large. In the absence of the testing data, it will not be possible to measure and detect overfitting. Overfitting is caused by generating a prediction function that has a seemingly perfect score on the samples but fails to predict on unseen data.*

- What does grid search do and why might you want to use it?

*Grid search is used to learn estimator parameters. For example: The max_depth parameter in the 'Decision Tree Regressor' can be optimized using an array of parameter values.*

- Why is cross validation useful and why might we use it with grid search?

*Cross-Validation is used to reduce the number of samples needed for learning the model. k-fold CV allows the training set to be split in 'k' smaller sets with (k-1) used for training and the last part for validation/testing. Further, the performance measure is looped for each of the k 'folds'.*

# 3) Analyzing Model Performance

- Look at all learning curve graphs provided. What is the general trend of training and testing error as training size increases?

*As training size increases training error rises and testing error falls.*

- Look at the learning curves for the decision tree regressor with max depth 1 and 10 (first and last learning curve graphs). When the model is fully trained does it suffer from either high bias/underfitting or high variance/overfitting?

  *Since the training error increases when the model is fully trained, my conclusion is that it suffers from high bias/underfitting.*

- Look at the model complexity graph. How do the training and test error relate to increasing model complexity? Based on this relationship, which model (max depth) best generalizes the dataset and why?

  *The training error reduces exponentially with increasing model complexity. The test error also reduces exponentially as the complexity increases from 0. However, it reaches a minimum value at max_depth = 6 and then fluctuates.*

  *max_depth = 6 seems to best generalize the dataset since the test error is the smallest.*

# 4) Model Prediction

- Model makes predicted housing price with detailed model parameters (max depth) reported using grid search. Note due to the small randomization of the code it is recommended to run the program several times to identify the most common/reasonable price/model complexity.

  *Most common max_depth = 4*

  *Prediction = 21.630*

- Compare prediction to earlier statistics and make a case if you think it is a valid model.

  *The value 21.630 is reasonable - It lies within the min/max range of [5.0 50.0].*

  *Also, it is within the standard deviation range (9.188) of the mean (22.533).*