

## Prediction of California Median House Prices Using Regression Techniques

### Import required libraries

```
In [42]: 1 %matplotlib inline
2 import seaborn as sns
3 import pandas as pd
4 import numpy as np
5 import os
6 import matplotlib.pyplot as plt
7 from sklearn import preprocessing
8 from sklearn.model_selection import train_test_split
9 from sklearn import tree
10 from sklearn import preprocessing
11 from sklearn.neighbors import KNeighborsClassifier
12 from sklearn.model_selection import train_test_split
13 from sklearn.metrics import precision_recall_fscore_support
14 from sklearn.metrics import accuracy_score
15 from sklearn import metrics
16 from sklearn.neighbors import KNeighborsClassifier
17 from sklearn import preprocessing
18 from sklearn.preprocessing import StandardScaler
19 from sklearn.linear_model import LogisticRegression, LogisticRegressionCV
20 from statsmodels.api as sm
21 from dmba import classificationSummary, gainChart, liftChart
22 from sklearn import metrics
23 from sklearn.model_selection import cross_val_score
24 from sklearn.model_selection import GridSearchCV
25 from sklearn.tree import DecisionTreeRegressor
26 from sklearn.linear_model import LogisticRegression
27 from sklearn.linear_model import DecisionTreeRegressor
28 from sklearn.metrics import mean_squared_error
29 from sklearn.preprocessing import MinMaxScaler
30 from sklearn.ensemble import RandomForestRegressor
31 from sklearn.linear_model import LinearRegression
32 import math
33 import gridfs
34 from scipy.stats import norm
35 from scipy.cluster import hierarchy as hc
36 from dmba import backward_elimination, forward_selection, stepwise_selection
37 from dmba import regressionSummary, exhaustive_search
38 from sklearn import linear_model
39 from sklearn.linear_model import Ridge
40 from regressionmetricraphplot import *
41 from sklearn.datasets import make_regression
42 from statsmodels.stats.outliers_influence import variance_inflation_factor
43 import warnings
44 warnings.filterwarnings('ignore')
```

### DATA SPLITTING (training, validation, and test sets)

```
In [43]: 1 df = pd.read_csv('data/California_Houses.csv')
2 df
```

```
Out[43]:
```

	Median_House_Value	Median_Income	Median_Age	Tot_Rooms	Tot_Bedrooms	Population	Households	Latitude	Longitude	Distance_to_coast	Distance
0	452600.0	8.252	41	880	129	322	126	37.88	-122.23	9263.040773	556529.
1	358500.0	8.3014	21	7099	1106	2401	1138	37.86	-122.22	10225.733072	554279.
2	352100.0	7.2574	52	1467	190	496	177	37.85	-122.24	8259.085109	554610.
3	341300.0	5.6431	52	1274	235	558	219	37.85	-122.25	7768.086571	555194.
4	342200.0	3.8462	52	1627	280	568	259	37.85	-122.25	7768.086571	555194.
...	...	...	...	...	...	...	...	...	...	...	...
20635	78100.0	1.5603	25	1665	374	845	330	39.48	-121.09	162031.481121	654530.
20636	77100.0	2.5568	18	697	150	356	114	39.49	-121.21	160445.433537	659747.
20637	92300.0	1.7000	17	2254	485	1007	433	39.43	-121.22	153754.341182	654042.
20638	84700.0	1.8672	18	1860	409	741	349	39.43	-121.32	152005.022339	657698.
20639	89400.0	2.3886	16	2785	616	1387	530	39.37	-121.24	146866.196892	648723.

20640 rows × 14 columns

After parsing through the initial data set, we can see there are 20640 entries and 14 data columns. All data represents a block within a neighborhood.

```
In [44]: 1 df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 14 columns):
 # Column          Non-Null Count  Dtype  
--- 
 0 Median_House_Value    20640 non-null   float64
 1 Median_Income        20640 non-null   float64
 2 Median_Age           20640 non-null   int64  
 3 Tot_Rooms            20640 non-null   int64  
 4 Tot_Bedrooms         20640 non-null   int64  
 5 Population           20640 non-null   float64
 6 Households           20640 non-null   int64  
 7 Latitude              20640 non-null   float64
 8 Longitude             20640 non-null   float64
 9 Distance_to_coast     20640 non-null   float64
 10 Distance_to_LA        20640 non-null   float64
 11 Distance_to_SanDiego  20640 non-null   float64
 12 Distance_to_SanJose   20640 non-null   float64
 13 Distance_to_SanFrancisco 20640 non-null   float64
dtypes: float64(9), int64(5)
memory usage: 2.2 MB
```

Initial descriptive statistics of the data. While not particularly useful now, it is important to see variances in the data.

```
In [45]: 1 df.describe()
```

```
Out[45]:
```

	Median_House_Value	Median_Income	Median_Age	Tot_Rooms	Tot_Bedrooms	Population	Households	Latitude	Longitude	Distance_to_coast	Distance
count	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000	20640.000000
mean	206855.816909	3.870671	28.639486	2635.763081	537.98014	1425.476744	499.539680	35.631861	-119.569704	40509.26	
std	115395.618674	1.899622	12.585558	2181.619282	421.247906	1132.462122	382.329753	2.139592	2.003532	49140.02	
min	14999.000000	0.49900	1.000000	2.000000	1.000000	3.000000	1.000000	32.540000	-124.350000	120.67	
25%	119600.000000	2.563400	18.000000	1447.750000	295.000000	787.000000	280.000000	33.930000	-121.800000	9079.75	
50%	179700.000000	3.534800	29.000000	2127.000000	435.000000	1166.000000	409.000000	34.260000	-118.490000	20522.01	
75%	264725.000000	4.74350	37.000000	3148.000000	647.000000	1725.000000	605.000000	37.710000	-118.010000	49830.41	
max	500001.000000	15.000100	52.000000	39320.000000	6445.000000	35682.000000	6082.000000	41.950000	-114.310000	333804.68	

Let's dive into the data and look for unique values to ensure there are no duplicates.

```
In [46]: 1 # Finding out how many unique values there are
2 dict = {}
3 for col in list(df.columns):
4     dict[col] = df[col].value_counts().shape[0]
5 pd.DataFrame(dict, index=['unique count']).T
```

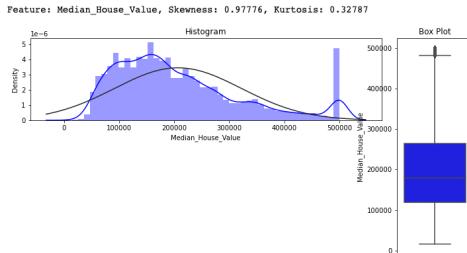
```
Out[46]:
```

	unique count
Median_House_Value	3842
Median_Income	12928
Median_Age	52
Tot_Rooms	5926
Tot_Bedrooms	1928
Population	3888
Households	1815
Latitude	862
Longitude	844
Distance_to_coast	12590
Distance_to_LA	12590
Distance_to_SanDiego	12590
Distance_to_SanJose	12590
Distance_to_SanFrancisco	12590

For initial pre-processing, let's group by data. Cont\_col is all data, Loc\_col is geographical and dist\_col is distance. The target variable is Median\_House\_Value.

```
In [47]: com_col = ['Median_Income', 'Median_Age', 'Tot_Rooms', 'Tot_Bedrooms', 'Population', 'Households', 'Distance_to_cities', 'Distance_to_LA', 'Distance_to_SanDiego', 'Distance_to_SanJose', 'Distance_to_SanFrancisco']
loc_col = ['Latitude', 'Longitude']
dist_col = ['Distance_to_Coast', 'Distance_to_LA', 'Distance_to_SanDiego', 'Distance_to_SanJose', 'Distance_to_SanFrancisco', 'Median_House_Value']
target = 'Median_House_Value'
```

```
In [48]: 1 def visualize_target(df, feature):
2     print('Feature: {} Skewness: {}, Kurtosis: {}'.format(feature,round(df[feature].skew(),5),round(df[feature].kurt(),5)))
3
4     fig = plt.figure(constrained_layout=True, figsize=(12,6))
5     grid = gridspec.GridSpec(ncols=5, nrows=5, figure=fig)
6
7     ax1 = fig.add_subplot(grid[0:2, :4])
8     ax1.set_title('Histogram')
9     sns.distplot(df[feature], hist=True, norm_hist=True, fit='norm', ax=ax1,color='blue')
10
11    ax2 = fig.add_subplot(grid[1:, :4])
12    ax2.set_title('Box Plot')
13    sns.boxplot(y=df[feature], orient='v', ax=ax2,color='blue')
14
15 visualize_target(df,target)
```



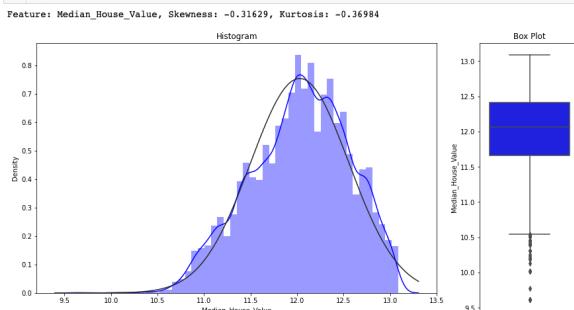
You can see the outliers on the higher end of the median house spectrum. We will now remove these according IQR.

```
In [49]: 1 def remove_outlier(df,col_name):
2     df[col_name].quantile(0.25)
3     q1 = df[col_name].quantile(0.75)
4     iqr = q3-q1 #Interquartile range
5     outlier = []
6     for i in df[col_name]:
7         if i <= (q1 - 1.5 * iqr) or i >=(q3 + 1.5 * iqr):
8             outlier.append(i)
9     outlier = pd.DataFrame(outlier,columns=['outlier'])
10    print("The Outliers removed in {} according to IQR :".format(round((outlier.shape[0])/df.shape[0])*100,2),col_name)
11    fence_low = q1-1.5*iqr
12    fence_high = q3+1.5*iqr
13    df = df[(df[col_name] > fence_low) & (df[col_name] < fence_high)]
14    return df
15
16 df = remove_outlier(df,target)
```

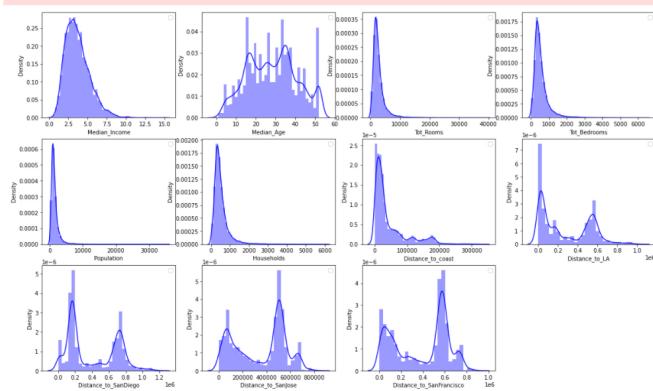
5.19% Outlier removed in Median\_House\_Value according to IQR

```
In [50]: df['Median_House_Value'] = np.log(df['Median_House_Value'])
```

This visualization is to check that all cuts



```
In [53]: 1 plt.figure(figsize=(10, 12))
2 for i, column in enumerate(cont_col, 1):
3     plt.subplot(3, 4, i)
4     sns.distplot(df[column], color="blue")
5     plt.legend()
6     plt.xlabel(column)
```



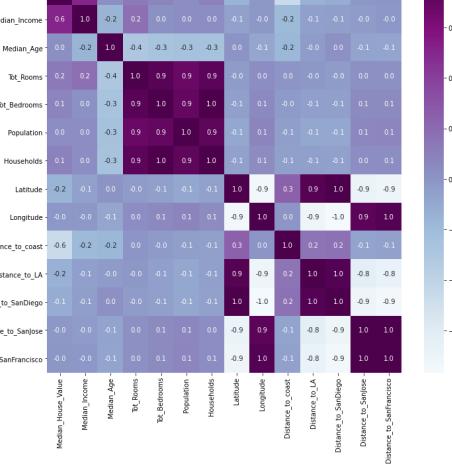
As you can see, all data columns have no extreme outliers. Exploratory Data Analysis can begin now.

## Exploratory Data Analysis

```
In [55]: 1 #Heatmap of a Correlation Table with new data frame
2 corr = df.corr()
3 fig, ax = plt.subplots()
4 fig.set_size_inches(12,11)
5 sns.heatmap(corr, annot=True, fmt=".1f", cmap="BuPu", center=0, ax=ax)
```

© 1995-2013

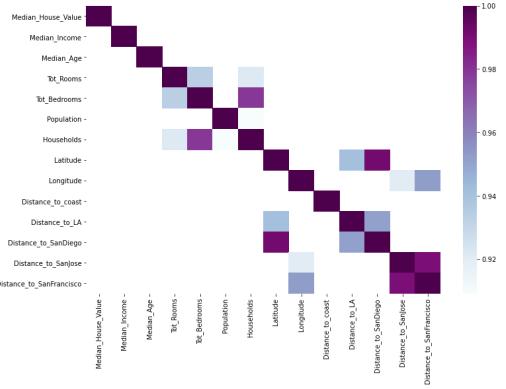




This first correlation is all data columns against each other. We are examining for all variables that are the highest correlation. Below we filter out all variables with greater than 0.9.

```
In [56]: 1 #Filtering out highly correlated variables (>0.9)with heatmap
2 corr = df.corr()
3 kpt = corr[corr>=.9]
4 plt.figure(figsize=(12,8))
5 sns.heatmap(kpt, cmap="BuPu")
```

Out[56]: <AxesSubplot:>



Highly correlated variables (>0.9) were observed among Tot\_Rooms, Tot\_Bedrooms, Households, and Population. Tot\_Room, Tot\_Bedrooms, Households will be excluded since Household variable alone can represent them.

Another set of high correlation between geographical variables: Latitude, Longitude, Distance\_to\_SanDiego, Distance\_to\_SanJose, and Distance\_to\_SanFrancisco. Thus, Latitude, Longitude, Distance\_to\_SanJose, and Distance\_to\_SanDiego will be dropped. Keeping only Distance\_to\_LA will cover Distance\_to\_SanDiego and Distance\_to\_SanFrancisco will cover Distance\_to\_SanJose since they are in close proximity.

```
In [57]: 1 #dropping highly correlated variables
2 df = df.drop(columns= ['Latitude', 'Longitude', 'Distance_to_SanJose', 'Tot_Bedrooms',
3 'Tot_Rooms', 'Population', 'Distance_to_SanDiego'])
4
5 df #new data frame after removing highly correlated attributes
```

Out[57]:

	Median_House_Value	Median_Income	Median_Age	Households	Distance_to_coast	Distance_to_LA	Distance_to_SanFrancisco
0	13.022764	8.3252	41	126	9263.040773	556629.158342	21260.137367
1	12.789684	8.3014	21	1138	10225.733072	554279.850069	20880.600400
2	12.771671	7.2574	52	177	8259.085109	554610.717069	18811.487450
3	12.740517	5.6431	52	219	7768.086571	555194.266086	18031.047568
4	12.743151	3.8462	52	259	7768.086571	555194.266086	18031.047568
...	...	...	...	...	...	...	...
20635	11.265745	1.5603	25	330	162031.481121	545530.186299	222619.890417
20636	11.252859	2.5568	18	114	160445.433537	559747.068444	21814.424634
20637	11.432799	1.7000	17	433	153754.341182	65042.214202	212097.986232
20638	11.346871	1.8672	18	349	152005.022239	657698.007703	207933.199166
20639	11.400876	2.3886	16	530	146866.196892	648723.337126	205473.376575

19569 rows × 7 columns

## DATA SPLITTING (training, validation, and test sets)

There were total records of 19569 and seven attributes after removing the outliers and highly correlated variables. The data set was then partitioned into training, validation, and test data sets (40%, 35%, and 25% respectively). There were 7828 records in training, 6849 records in validation, and 4620 records in test data sets with 7 variables.

```
In [58]: 1 #Total records after removing outliers = 19569,
2 #Splitting train = 7828(40%), valid = 6849(35%), test = 4620(25%)
3
4 train, temp = train_test_split(df, train_size= 7828, random_state=1)
5 valid, test = train_test_split(temp, train_size= 6849, random_state=1)
6 print("Training : ", train.shape)
7 print("Validation : ", valid.shape)
8 print("Test : ", test.shape)
```

Training : (7828, 7)  
Validation : (6849, 7)  
Test : (4620, 7)

Model will be trained on the train data set and evaluated with validation data set for better performance. The best model will be selected by comparing the metrics of each model on the test data set.

The selected models were built to predict the target outcome of Median House Value by using the independent predictors: 'Median\_Income', 'Median\_Age', 'Households', 'Distance\_to\_coast', 'Distance\_to\_LA', 'Distance\_to\_SanFrancisco'. The target outcome was the Median House Value.

```
In [59]: 1 #Assigning predictors and outcome variables
2 predictors = list(df.columns)
3 outcome = 'Median_House_Value'
4 predictors.remove(outcome)
5 print(predictors)
6 print(outcome)
7
8 #Assigning Xs,y for train, valid, and test data sets
9 train_X = train[predictors] #train_X, train_y for train data set
10 train_y = train[outcome]
11 valid_X = valid[predictors] #valid_X, valid_y for valid data set
12 valid_y = valid[outcome]
13 X = test[predictors] #X, y for test data set
```

```

14 y = test[outcome]
['Median_Income', 'Median_Age', 'Households', 'Distance_to_coast', 'Distance_to_LA', 'Distance_to_SanFrancisco']
Median_House_Value

```

## MULTIPLE LINEAR REGRESSION MODEL

For the Multiple Linear Regression model, the chosen parameter of stepwise regression was used for the best predictors.

```

In [60]: 1 #Multiple Linear Regression Model with stepwise regression for selecting the best variables
2
3 def train_model(variables):
4     if len(variables) == 0:
5         return None
6     model = LinearRegression()
7     return model.fit(train_X[variables], train_y)
8
9 def score_model(model, variables):
10    if len(variables) == 0:
11        return mean_squared_error(train_y, [train_y.mean()] * len(train_y), model)
12    return mean_squared_error(train_y, model.predict(train_X[variables]))
13
14
15 linear_model, best_variables1 = stepwise_selection(predictors, train_model, score_model,
16                                                    direction='stepwise', verbose=True)
17 print(best_variables1)
18 linear_predictors = best_variables1
19
Variables: Median_Income, Median_Age, Households, Distance_to_coast, Distance_to_LA, Distance_to_SanFrancisco
Start: score=0.28, constant
Step: score=0.16, add Median_Income
Step: score=0.11, add Distance_to_coast
Step: score=0.11, add Households
Step: score=0.11, add Median_Age
Step: score=0.11, add Distance_to_LA
Step: score=0.11, add Distance_to_SanFrancisco
Step: score=0.11, unchanged None
['Median_Income', 'Distance_to_coast', 'Households', 'Median_Age', 'Distance_to_LA', 'Distance_to_SanFrancisco']

The best predictor was Median_Income with a score of 01.6, and the rest were Distance_to_coast, Households, Median_Age, Distance_to_LA, Distance_to_San_Francisco with the same score of 0.11. All the predictors were used for training the model since all revealed similar scores.

```

```

In [61]: 1 #Regression Statistics by MULTIPLE LINEAR REGRESSION MODEL
2 print('Train Performance')
3 regressionSummary(train_y, linear_model.predict(train_X[best_variables1]))
4
5 print('Validation Performance')
6 regressionSummary(valid_y, linear_model.predict(valid_X[best_variables1]))
7
8 print('Test Performance')
9 regressionSummary(y, linear_model.predict(X[best_variables1]))

```

Train Performance

Regression statistics

```

Mean Error (ME) : -0.0000
Root Mean Squared Error (RMSE) : 0.3265
Mean Absolute Error (MAE) : 0.1614
Mean Percentage Error (MPE) : -0.0750
Mean Absolute Percentage Error (MAPE) : 2.1019
Validation Performance

```

Regression statistics

```

Mean Error (ME) : -0.0043
Root Mean Squared Error (RMSE) : 0.3250
Mean Absolute Error (MAE) : 0.2511
Mean Percentage Error (MPE) : -0.1144
Mean Absolute Percentage Error (MAPE) : 2.1027
Test Performance

```

Regression statistics

```

Mean Error (ME) : 0.0045
Root Mean Squared Error (RMSE) : 0.3183
Mean Absolute Error (MAE) : 0.2474
Mean Percentage Error (MPE) : -0.0365
Mean Absolute Percentage Error (MAPE) : 2.0652

```

## REGRESSION TREE MODEL

The Stepwise regression was also used for the Regression Tree model.

```

In [62]: 1 # regression tree
2 def train_model(variables):
3     if len(variables) == 0:
4         return None
5     model = DecisionTreeRegressor()
6     return model.fit(train_X[variables], train_y)
7
8 def score_model(model, variables):
9     if len(variables) == 0:
10        return mean_squared_error(train_y, [train_y.mean()] * len(train_y), model)
11    return mean_squared_error(train_y, model.predict(train_X[variables]))
12
13 tree_model, best_variables = stepwise_selection(predictors, train_model, score_model, direction='forward', verbose=True)
14 print(best_variables)

Variables: Median_Income, Median_Age, Households, Distance_to_coast, Distance_to_LA, Distance_to_SanFrancisco
Start: score=0.28, constant
Step: score=0.01, add Distance_to_coast
Step: score=0.00, add Median_Income
Step: score=0.00, unchanged None
['Distance_to_coast', 'Median_Income']

The model was built using the best predictors, Distance_to_coast and Median_Income.

```

```

In [63]: 1 #Regression Statistics by REGRESSION TREE MODEL
2 print('Train Performance')
3 regressionSummary(train_y, tree_model.predict(train_X[best_variables]))
4 print('Validation performance')
5 regressionSummary(valid_y, tree_model.predict(valid_X[best_variables]))
6 print('Test performance')
7 regressionSummary(y, tree_model.predict(X[best_variables]))

```

Train Performance

Regression statistics

```

Mean Error (ME) : 0.0000
Root Mean Squared Error (RMSE) : 0.0000
Mean Absolute Error (MAE) : 0.0000
Mean Percentage Error (MPE) : 0.0000
Mean Absolute Percentage Error (MAPE) : 0.0000
Validation performance

```

Regression statistics

```

Mean Error (ME) : -0.0180
Root Mean Squared Error (RMSE) : 0.4373
Mean Absolute Error (MAE) : 0.3285
Mean Percentage Error (MPE) : -0.0208
Mean Absolute Percentage Error (MAPE) : 2.7479
Test performance

```

Regression statistics

```

Mean Error (ME) : -0.0049
Root Mean Squared Error (RMSE) : 0.4417
Mean Absolute Error (MAE) : 0.3343
Mean Percentage Error (MPE) : -0.1088
Mean Absolute Percentage Error (MAPE) : 2.7865

```

Since it was all zero percent error for train data set, overfitting was noted for Regression Tree model.

## RANDOM FOREST

Random Forest model was built using n\_estimators of 500 and random\_state 1. The two most important features were Distance\_to\_coast and Median\_Income.

```

In [64]: 1 rf = RandomForestRegressor(n_estimators=500, random_state=1)
2 rf.fit(train_X, train_y)

```

```
Out[64]: RandomForestRegressor(n_estimators=500, random_state=1)
```

```

In [65]: 1 #To see the most important features
2
3 importances = rf.feature_importances_
4 std = np.std([tree.feature_importances_ for tree in rf.estimators_], axis=0)
5
6 df = pd.DataFrame({'feature': train_X.columns, 'importance': importances, 'std': std})
7 df = df.sort_values('importance')
8 print(df)
9

```

```

10 ax = df.plot(kind='barh', xerr='std', x='feature', legend=False)
11 ax.set_ylabel('')
12 plt.tight_layout()
13 plt.show()

      feature importance std
1   Median_Age  0.034889  0.003350
2 Households  0.047419  0.004691
4 Distance_to_LA  0.091123  0.006367
5 Distance_to_SanFrancisco  0.034123  0.004876
0 Median_Income  0.343972  0.041653
3 Distance_to_coast  0.382924  0.042084

```

Random Forest model was built using n\_estimators of 500 and random\_state 1. The two most important features were Distance\_to\_coast and Median\_Income.

```

In [66]: 1 #Regression Statistics by RANDOM FOREST MODEL
2 regressionSummary(train_y, rf.predict(train_X))
3 regressionSummary(valid_y, rf.predict(valid_X))
4 regressionSummary(y, rf.predict(X))

```

```

Regression statistics
Mean Error (ME) : 0.0006
Root Mean Squared Error (RMSE) : 0.0889
Mean Absolute Error (MAE) : 0.0633
Mean Percentage Error (MPE) : -0.0057
Mean Absolute Percentage Error (MAPE) : 0.5281

Regression statistics
Mean Error (ME) : -0.0071
Root Mean Squared Error (RMSE) : 0.2435
Mean Absolute Error (MAE) : 0.1716
Mean Percentage Error (MPE) : -0.1043
Mean Absolute Percentage Error (MAPE) : 1.4368

Regression statistics

```

## Lasso Regression

We're going to test Ridge and Lasso regression as they're two simple techniques that help reduce model complexity and prevent over-fitting that might occur from simple linear regression.

```

In [67]: 1 lasso = Lasso()
2 lasso.fit(train_X,train_y)
3 train_score=lasso.score(train_X, train_y)
4 test_score=lasso.score(valid_X, valid_y)
5 coeff_used = np.sum(lasso.coef_!=0)
6
7 print ("training score:", train_score)
8 print ("test score:", test_score)
9 print ("number of features used:", coeff_used)

training score: 0.3544637691778752
test score: 0.3807762875569956
number of features used: 4

```

```

In [68]: 1 # we'll tune lasso to alpha = 0.01
2
3 lasso_001 = Lasso(alpha=0.01, max_iter=10e5)
4 lasso_001.fit(train_X, train_y)
5
6 train_score001=lasso_001.score(train_X, train_y)
7 test_score001=lasso_001.score(valid_X, valid_y)
8 coeff_used001 = np.sum(lasso_001.coef_!=0)
9
10 print ("training score for alpha=0.01", train_score001)
11 print ("test score for alpha=0.01", test_score001)
12 print ("number of features used for alpha =0.01:", coeff_used001)

training score for alpha=0.01: 0.616067695362032
test score for alpha =0.01: 0.6307185175715928
number of features used for alpha =0.01: 6

```

```

In [69]: 1 # we'll tune lasso to alpha = 0.0001
2
3 lasso_00001 = Lasso(alpha=0.0001, max_iter=10e5)
4 lasso_00001.fit(train_X, train_y)
5
6 train_score00001=lasso_00001.score(train_X, train_y)
7 test_score00001=lasso_00001.score(valid_X, valid_y)
8 coeff_used00001 = np.sum(lasso_00001.coef_!=0)
9
10 print ("training score for alpha=0.0001", train_score00001)
11 print ("test score for alpha =0.0001", test_score00001)
12 print ("number of features used for alpha =0.0001:", coeff_used00001)

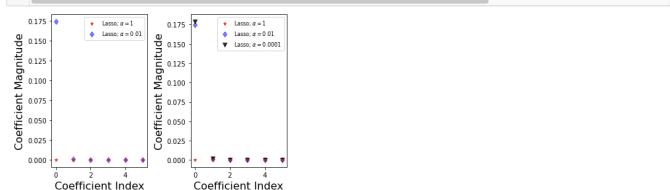
training score for alpha=0.0001: 0.416250494217306
test score for alpha =0.0001: 0.6312128947118348
number of features used for alpha =0.0001: 6

```

```

In [70]: 1 # plot the coefficients
2
3 plt.subplot(1,2,1)
4 plt.plot(lasso.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,color='red',label=r'Lasso; $\alpha$ = 1',zorder=1)
5 plt.plot(lasso_001.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue',label=r'Lasso; $\alpha$ = 0.01',zorder=2)
6 plt.xlabel('Coefficient Index',fontsize=16)
7 plt.ylabel('Coefficient Magnitude',fontsize=16)
8 plt.legend(fontsize=16)
9 plt.subplot(1,2,2)
10 plt.plot(lasso.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,color='red',label=r'Lasso; $\alpha$ = 1',zorder=1)
11 plt.plot(lasso_00001.coef_,alpha=0.5,linestyle='none',marker='d',markersize=6,color='blue',label=r'Lasso; $\alpha$ = 0.0001',zorder=2)
12 plt.plot(lasso_0001.coef_,alpha=0.8,linestyle='none',marker='v',markersize=6,color='black',label=r'Lasso; $\alpha$ = 0.001',zorder=3)
13 plt.xlabel('Coefficient Index',fontsize=16)
14 plt.ylabel('Coefficient Magnitude',fontsize=16)
15 plt.legend(fontsize=16)
16 plt.tight_layout()
17 plt.show()

```



```

In [71]: 1 #Regression Statistics by Lasso Regression MODEL
2 regressionSummary(train_y, lasso_00001.predict(train_X))
3 regressionSummary(valid_y, lasso_00001.predict(valid_X))
4 regressionSummary(y, lasso_00001.predict(X))

```

```

Regression statistics
Mean Error (ME) : 0.0000
Root Mean Squared Error (RMSE) : 0.3265
Mean Absolute Error (MAE) : 0.2516
Mean Percentage Error (MPE) : -0.0750
Mean Absolute Percentage Error (MAPE) : 2.1019

Regression statistics
Mean Error (ME) : -0.0043
Root Mean Squared Error (RMSE) : 0.3250
Mean Absolute Error (MAE) : 0.2511
Mean Percentage Error (MPE) : -0.1144
Mean Absolute Percentage Error (MAPE) : 2.1027

Regression statistics

```

```

Mean Error (ME) : 0.0045
Root Mean Squared Error (RMSE) : 0.3183
Mean Absolute Error (MAE) : 0.2474
Mean Percentage Error (MPE) : -0.0365
Mean Absolute Percentage Error (MAPE) : 2.0652

```

### Ridge Regression

In [72]:

```

1 # here we'll tune alpha to 0.01 & 100, the higher the alpha, more limitations on the coefficients.
2
3 # alpha = 0.1
4 RR = Ridge(alpha=0.01)
5 RR.fit(train_X, train_y)
6
7 Ridge.train_score=RR.score(train_X, train_y)
8 Ridge.test_score=RR.score(valid_X, valid_y)
9 coeff_used = np.sum(lasso.coef_!=0)
10
11 print ("training score:", Ridge.train_score)
12 print ("test score: ", Ridge.test_score)
13 print ("number of features used: ", coeff_used)

training score: 0.6162805177038775
test score: 0.631216207742185
number of features used: 4

```

In [73]:

```

1 # alpha = 100
2 rr100 = Ridge(alpha=100) # comparison with alpha value
3 rr100.fit(train_X, train_y)
4
5 Ridge.train_score100 = rr100.score(train_X, train_y)
6 Ridge.test_score100 = rr100.score(valid_X, valid_y)
7 coeff_used = np.sum(lasso.coef_!=0)
8
9 print ("training score:", Ridge.train_score100)
10 print ("test score: ", Ridge.test_score100)
11 print ("number of features used: ", coeff_used)

training score: 0.6162419189776894
test score: 0.6311039651446302
number of features used: 4

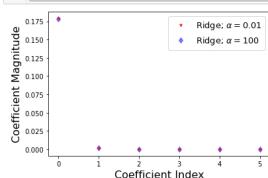
```

In [74]:

```

1 # plot the coefficients
2
3 plt.plot(RR.coef_,alpha=0.1,linestyle='none',marker='*',markersize=5,color='red',label='Ridge; $\alpha = 0.01$',
4          picker=rr100.coef_.alpha+0.5,linestyle='none',marker='d',markersize=6,color='blue',label='Ridge; $\alpha = 100$')
5 plt.xlabel('Coefficient Index',fontsize=16)
6 plt.ylabel('Coefficient Magnitude',fontsize=16)
7 plt.legend(fontsize=13,loc=1)
8 plt.show()

```



In [75]:

```

1 #Regression Statistics by Ridge Regression MODEL
2 regressionSummary(train_y, rr100.predict(train_X))
3 regressionSummary(valid_y, rr100.predict(valid_X))
4 regressionSummary(y, rr100.predict(X))

```

Regression statistics

```

Mean Error (ME) : 0.0000
Root Mean Squared Error (RMSE) : 0.3285
Mean Absolute Error (MAE) : 0.2516
Mean Percentage Error (MPE) : -0.0752
Mean Absolute Percentage Error (MAPE) : 2.1019

```

Regression statistics

```

Mean Error (ME) : -0.0043
Root Mean Squared Error (RMSE) : 0.3250
Mean Absolute Error (MAE) : 0.2511
Mean Percentage Error (MPE) : -0.1130
Mean Absolute Percentage Error (MAPE) : 2.1030

```

Regression statistics

```

Mean Error (ME) : 0.0045
Root Mean Squared Error (RMSE) : 0.3183
Mean Absolute Error (MAE) : 0.2474
Mean Percentage Error (MPE) : -0.0369
Mean Absolute Percentage Error (MAPE) : 2.0653

```

## RESULTS AND FINAL MODEL SELECTION

In [76]:

```

1 #REGRESSION STATISTICS OF DIFFERENT MODELS (BASED ON THE TEST DATA)
2
3 print("MULTIPLE LINEAR REGRESSION:")
4 regressionSummary(y,linear_model.predict(X[best_variables]))
5
6 print("LASSO:")
7 regressionSummary(y,lasso_00001.predict(X))
8
9 print("RIDGE:")
10 regressionSummary(y,rr100.predict(X))
11
12 print("REGRESSION TREE MODEL:")
13 regressionSummary(y, tree_model.predict(X[best_variables]))
14
15 print("RANDOM FOREST:")
16 regressionSummary(y, rf.predict(X))
17
18 #According to Regression Statistics and Residual Errors' plot, RANDOM FOREST DEEMS THE BEST MODEL.

```

MULTIPLE LINEAR REGRESSION:

Regression statistics

```

Mean Error (ME) : 0.0045
Root Mean Squared Error (RMSE) : 0.3183
Mean Absolute Error (MAE) : 0.2474
Mean Percentage Error (MPE) : -0.0365
Mean Absolute Percentage Error (MAPE) : 2.0652
LASSO:

```

Regression statistics

```

Mean Error (ME) : 0.0045
Root Mean Squared Error (RMSE) : 0.3183
Mean Absolute Error (MAE) : 0.2474
Mean Percentage Error (MPE) : -0.0365
Mean Absolute Percentage Error (MAPE) : 2.0652
RIDGE:

```

Regression statistics

```

Mean Error (ME) : 0.0045
Root Mean Squared Error (RMSE) : 0.3183
Mean Absolute Error (MAE) : 0.2474
Mean Percentage Error (MPE) : -0.0369
Mean Absolute Percentage Error (MAPE) : 2.0653
REGRESSION TREE MODEL:

```

Regression statistics

```

Mean Error (ME) : -0.0049
Root Mean Squared Error (RMSE) : 0.4417
Mean Absolute Error (MAE) : 0.3343
Mean Percentage Error (MPE) : -0.1088
Mean Absolute Percentage Error (MAPE) : 2.7865
RANDOM FOREST:

```

Regression statistics

```

Mean Error (ME) : 0.0007
Root Mean Squared Error (RMSE) : 0.2330
Mean Absolute Error (MAE) : 0.1647
Mean Percentage Error (MPE) : -0.0341
Mean Absolute Percentage Error (MAPE) : 1.3736

```

In [77]:

```

1 #Printing R2 and RMSE metrics from each model
2
3 pred_test = linear_model.predict(X[best_variables])
4 pred_RT_test = tree_model.predict(X[best_variables])
5 pred_RF_test = rf.predict(X)
6
7 print('Multiple Linear Regression:', CompareModels.R2AndRMSE(y_true=y, y_pred=pred_test))
8

```

```

9 print('Random Forest:', CompareModels.R2AndRMSE(y_test=y, y_pred= rf.predict(X)))
10 print('Regression Tree:', CompareModels.R2AndRMSE(y_test=y, y_pred=pred_RT_test))
11 print('Lasso Regression:', CompareModels.R2AndRMSE(y_test=y, y_pred = lasso_00001.predict(X)))
12 print ('Ridge Regression:', CompareModels.R2AndRMSE(y_test =y, y_pred = rr100.predict(X)))
13
14
15
16
17
```

Multiple Linear Regression: (0.6282688626237674, 0.31829977423969413)  
Random Forest: (0.805854849945904, 0.2302034691541602)  
Regression Tree: (0.2843100370921403, 0.4416561939597983)  
Lasso Regression: (0.628266086798013, 0.3183073922311347)  
Ridge Regression: (0.628227390382614, 0.3183175293261262)

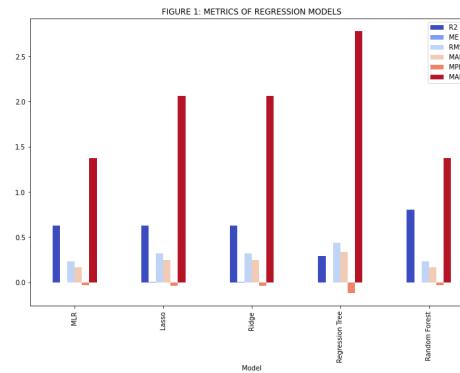
```

In [78]: 1 #Creating Summary Metrics Table of each model on test data set
2
3 data = {'Model': ['MLR', 'Lasso', 'Ridge', 'Regression Tree', 'Random Forest'],
4         'R2': [0.6283, 0.6283, 0.6282, 0.2886, 0.8059],
5         'ME': [0.0007, 0.0045, 0.0045, -0.0059, 0.0007],
6         'RMSE': [0.2300, 0.3183, 0.3183, 0.4403, 0.2300],
7         'MAE': [-0.0344, 0.0318, 0.0318, 0.0747, -0.0344],
8         'MPE': [-0.0344, 0.0318, 0.0318, -0.1172, -0.0344],
9         'MAPE': [1.3736, 2.0652, 2.0653, 2.7772, 1.3736],
10        }
11 Metric_Table = pd.DataFrame(data)
12 Metric_Table
```

	Model	R2	ME	RMSE	MAE	MPE	MAPE
0	MLR	0.6283	0.0007	0.2300	0.1647	-0.0341	1.3736
1	Lasso	0.6283	0.0045	0.3183	0.2474	-0.0365	2.0652
2	Ridge	0.6282	0.0045	0.3183	0.2474	-0.0369	2.0653
3	Regression Tree	0.2886	-0.0059	0.4403	0.3384	-0.1172	2.7772
4	Random Forest	0.8059	0.0007	0.2300	0.1647	-0.0341	1.3736

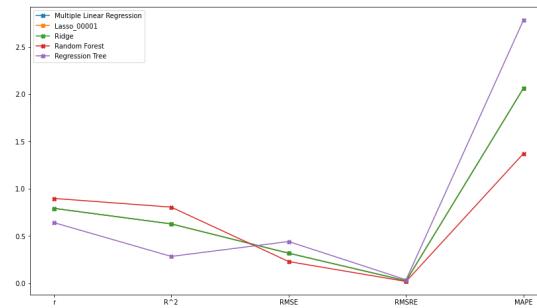
```

In [79]: 1 #Bargraph for Comparison of Metrics for Each Model
2
3 Metric_Table.plot(x="Model", y=["R2", "ME", "RMSE", "MAE", "MPE", "MAPE"], kind="bar", figsize=(12,8), colormap='coolwarm')
4 plt.title("FIGURE 1: METRICS OF REGRESSION MODELS")
5 plt.show()
```



```

In [80]: 1 #Plotting graphs to show the best model and similarity between MLR, LASSO, and RIDGE.
2
3 plot = CompareModels()
4
5 plot.add(model_name = 'Multiple Linear Regression', y_test=y, y_pred= pred_test)
6 plot.add(model_name = 'Lasso_00001', y_test =y, y_pred = lasso_00001.predict(X))
7 plot.add(model_name = 'Ridge', y_test =y, y_pred = rr100.predict(X))
8
9 plot.add(model_name='Random Forest', y_test=y, y_pred=pred_RF_test)
10
11 plot.add(model_name='Regression Tree', y_test=y, y_pred=pred_RT_test)
12
13 plot.add(model_name='Random Forest', y_test=y, y_pred=pred_RF_test)
14
15
16 plot.show(figsize=(14,8))
17
18 #RANDOM FOREST PERFORMED THE BEST.
19 #THE METRICS FOR MLR, RIDGE, AND LASSO WERE ALMOST IDENTICAL AND OVERLAPPED IN THE GRAPH.
20 #REGRESSION TREE YIELDED THE LEAST R2 VALUE AND HIGHEST ERROR RATES.
```



From the Metric Table, metrics of Multiple Linear Regression, Ridge, and Lasso were almost identical with R2 value of 0.6282 (62.82%). Thus, in the above graph, all three overlap together and doesn't show any difference between these three models.

The metrics of Regression Tree for the train data set showed zero percent error revealing the overfitting behavior and the R2 value of 28.86%, the lowest score out of all the models.

Among them, Random Forest scored the highest R2 value of 80.59% best fitting the data set. The metrics for the Random Forest also ranked the lowest percentage errors of ME, RMSE, MAE, MPE, and MAPE. Therefore, we would use this algorithm for all machine learning processes. The application of this research can be applied to applications across the real estate industry. The models had high precision outputs and could be utilized to create deliverables that can identify real estate evaluation to answer questions such as: Which house should I buy or build to maximize my return? Where or when should I do so? What is its optimum rent or sale price?

## Recommended Next Steps

The Random Forest scored the highest R2 value of 80.59% best fitting the data set. The metrics for the Random Forest also ranked the lowest percentage errors of ME, RMSE, MAE, MPE, and MAPE. Therefore, we would use this algorithm for all machine learning processes. The application of this research can be applied to applications across the real estate industry. The models had high precision outputs and could be utilized to create deliverables that can identify real estate evaluation to answer questions such as: Which house should I buy or build to maximize my return? Where or when should I do so? What is its optimum rent or sale price?

In [ ]: