

1 Start coding or [generate](#) with AI.

Assignment Supervised machine Learning - Regression

House Price Prediction Using Bolton Housing Data

This dataset provides information about house prices in Bolton. House Price Prediction

In this notebook, you'll follow the basic machine learning process to build a regression model to predict house prices using the "Boston Housing Dataset" from sklearn. The regression model will either be a Decision Tree or Random Forest regressor.

Follow the instructions and complete each TODO to complete the assessment on the essential steps in building and evaluating a regression model.

The following is a description of each column in the dataset:

Dataset Features (Bolton Housing):

- CRIM: Crime rate by town
- ZN: Proportion of residential land zoned for large lots
- INDUS: Proportion of non-retail business acres per town
- CHAS: Charles River dummy variable (1 if tract bounds river; 0 otherwise)
- NOX: Nitric oxide concentration (parts per 10 million)
- RM: Average number of rooms per dwelling
- AGE: Proportion of owner-occupied units built before 1940
- DIS: Weighted distances to five Boston employment centers
- RAD: Index of accessibility to radial highways
- TAX: Full-value property tax rate per \$10,000
- PTRATIO: Pupil-teacher ratio by town
- B: Proportion of Black population
- LSTAT: Percentage of lower status of the population
- MEDV (Target): Median value of owner-occupied homes in \$1,000s

Dataset is from sklearn Datasets

```

1 from seaborn.palettes import color_palette
2
3
4 # --- Imports ---
5 # TODO: Import all the necessary libraries for data handling, visualization, and model building.
6 # Example: import pandas as pd
7 # Add your imports here:
8
9 import pandas as pd
10 import numpy as np
11 from sklearn.datasets import fetch_openml
12 from sklearn.model_selection import train_test_split
13 from sklearn.tree import DecisionTreeRegressor
14 from sklearn.ensemble import RandomForestRegressor
15 from sklearn.metrics import mean_squared_error, r2_score
16 from matplotlib import pyplot as plt
17 import seaborn as sns
18 from sklearn.model_selection import GridSearchCV
19
20 # set up plot area:
21
22 # %matplotlib inline
23 # sns.set_style(style='darkgrid', palette='pastel')
24
25 # --- Data Collection and Loading ---
26 # TODO: Load the 'Boston Housing' dataset from sklearn and convert it into a pandas DataFrame.
27 # Hint: Use `load_boston()` from `sklearn.datasets`
28
29 # Load dataset and convert to DataFrame:
30
31 # Add your code here:
32

```

```

33 boston_data = fetch_openml(name='boston', version=1, as_frame=True) # load the boston dataset from sklearn
34 df = boston_data.frame # convert to df
35
36 # --- Quick Check of Data ---
37 # TODO: Display the first few rows of the dataset to understand its structure.
38 # Hint: Use `.head()` to inspect the first few rows.
39
40 # Add your code here:
41 print("Inspect the Dataset")
42 display(df.head())
43 print("-" * 30)
44 # --- Data Exploration ---
45
46 # TODO: Check the features and target variable. Identify which is continuous and categorical if applicable.
47 # Hint: Use `.info()` and `.describe()` to inspect data types and statistical properties.
48
49 # Add your code here:
50 print("Dataset Info:")
51 df.info()
52 print("Summary Statistics:")
53 print("-" * 140)
54 display(df.describe())
55 print("-" * 15)
56 # --- EDA and Data Preprocessing ---
57 # TODO: Check for missing/null values.
58 # Hint: Use `.isnull().sum()` to check for null values.
59
60 # Add your code here:
61
62 print("Missing Values:")
63 print(df.isnull().sum())
64 print("-" * 25)
65
66 # Elimiate Duplicates
67
68 df=df.drop_duplicates()
69 print("Duplicates Removed")
70 print("-" * 25)
71
72 # --- Step 5: Rename Columns for Clarity ---
73
74 # Rename columns to descriptive
75
76
77 df.rename(columns={
78     'CRIM': 'crime_rate',          # per capita crime rate by town
79     'ZN': 'residential_land_zone', # proportion of residential land zoned for lots over 25,000 sq.ft.
80     'INDUS': 'non_retail_business_acres', # proportion of non-retail business acres per town
81     'CHAS': 'charles_river_dummy',  # Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
82     'NOX': 'nitric_oxides_conc',    # nitric oxides concentration (parts per 10 million)
83     'RM': 'avg_rooms_per_dwelling', # average number of rooms per dwelling
84     'AGE': 'owner_occupied_units_pct', # proportion of owner-occupied units built prior to 1940
85     'DIS': 'distance_to_employment', # weighted distances to five Boston employment centers
86     'RAD': 'highway_accessibility_index', # index of accessibility to radial highways
87     'TAX': 'property_tax_rate',     # full-value property-tax rate per $10,000
88     'PTRATIO': 'pupil_teacher_ratio', # pupil-teacher ratio by town
89     'B': 'black_population_index',  # 1000(Bk - 0.63)^2 where Bk is proportion of Black population by town
90     'LSTAT': 'lower_status_pct',    # percentage of lower status population
91     'MEDV': 'median_home_value'     # median value of owner-occupied homes in $1000s
92 }, inplace=True)
93
94 # --- Display the renamed columns ---
95 print("Renamed Columns:")
96 print(df.columns.tolist())
97
98
99 # --- Data Visualization ---
100
101 # TODO: Visualize the data. Create scatter plots to see the relationship between independent features and the target variable.
102 # Example: Use `plt.scatter()` to visualize the relationship between features like 'RM' (average number of rooms) and the target variable 'MEDV'.
103
104 # Add your code here:
105 plt.figure(figsize=(10, 6))
106 sns.heatmap(
107     df.corr(),
108     annot=True, # indicates the corr values
109     cmap='coolwarm',
110     cbar=True,

```

```

110     tmt = .4T ,
111     linewidths=0.5, # add spaces between cells
112     cbar_kws={'shrink': 0.8}) #smaller color bar
113 plt.title('Correlation Heatmap', fontsize=14, pad=12)
114 plt.tight_layout()
115 plt.xticks(rotation=45) # rotate for readability
116 plt.yticks(rotation=0)
117 plt.show()
118
119
120 # TODO: Create a function to automate scatter plots for all features vs the target variable.
121 # Hint: The function should loop over a list of features and plot scatter plots for each.
122
123 # Define your function here:
124
125 def plot_scatter_matrix(df, target_col, feature_cols): #plotting the target col (median_home_value) against all independent
126     feature_cols = [c for c in df.select_dtypes(include=[np.number]).columns if c !=target_col] # get all the numerical col
127     plt.figure(figsize=(15, 12)) # set the figure size
128     plt.suptitle('Scatter Plots: Features vs Target (Median Home Value)', fontsize=16, y=1.02) # plot title- using this sup
129
130     for i, feature in enumerate(feature_cols, 1): # Loop through each feature - this code goes through each feature in the
131         plt.subplot(3, 5, i) # adjust the grid, 2 cols, 3 rows based on the features in the dataset
132         sns.scatterplot(x=df[feature], y=df[target_col], data=df, alpha=0.6, palette='teal')
133         plt.title(f'{feature_cols} vs {target_col}', fontsize=11)
134         plt.xlabel(feature_cols, fontsize=9) # plotting the x-axis
135         plt.ylabel(target_col, fontsize=9) #plotting the y-axis
136
137
138 plt.tight_layout()
139 plt.show()
140
141 # call the function
142
143 plot_scatter_matrix(df, target_col='median_home_value', feature_cols='feature')
144
145 # --- Feature Selection ---
146 # TODO: Use the function to visualize the relationships between multiple features and the target variable.
147 # Example: ['RM', 'LSTAT', 'AGE', 'CRIM']
148 # Target: 'MEDV'
149
150 # Add your code here:
151
152 selected_feature = ['crime_rate', 'residential_land_zone', 'non_retail_business_acres', 'nitric_oxides_conc']
153
154 def viz_selected_feat(df, features, target_col='median_home_value'):
155     plt.figure(figsize=(15, 12))
156     plt.suptitle('Selected Features vs Target (Median Home Value)', fontsize=16, y=1.02)
157
158     for j, feature in enumerate(features, 1):
159         plt.subplot(2, 2, j)
160         plt.scatter(x=df[feature], y=df[target_col], alpha=0.6, color='teal')
161         plt.title(f'{feature} vs {target_col}', fontsize=14, pad=12)
162         plt.xlabel(features, fontsize=9)
163         plt.ylabel(target_col, fontsize=9)
164
165 plt.tight_layout()
166 plt.show()
167
168 viz_selected_feat(df, target_col='median_home_value', features=selected_feature)
169 print("-" * 120)
170
171 # --- ML Model Training ---
172 # TODO: Split the dataset into training and testing sets.
173 # Hint: Use `train_test_split()` from `sklearn.model_selection` with an 80/20 split.
174
175 # Define X (features) and y (target) and perform the train-test split:
176
177 # Define feature matrix X (all independent variables)
178 # and target vector y (dependent variable)
179 X = df.drop('median_home_value', axis=1) # the target variable has been dropped, because it is the variable that I wish to p
180 y = df['median_home_value']
181
182 # split the dataset for training:
183
184 X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2, random_state=42)
185
186 # TODO: Choose an appropriate regression model: Decision Tree or Random Forest.
187 # Hint: Use either `DecisionTreeRegressor` or `RandomForestRegressor` from `sklearn.tree` or `sklearn.ensemble`.

```

```

188
189 # Define your regression model here:
190
191 model = DecisionTreeRegressor(random_state=42)
192
193 # TODO: Train the model on the training data.
194 # Hint: Use `.fit()` to train the model.
195
196 # Add your code here:
197
198 model.fit(X_train, y_train)
199
200 # --- Model Evaluation ---
201 # TODO: Evaluate the performance of the model on the test set using relevant metrics (e.g., RMSE, R-squared).
202 # Hint: Use `mean_squared_error()` and `r2_score()` from `sklearn.metrics`.
203
204 # Predict on the test set and calculate the evaluation metrics:
205
206 y_pred = model.predict(X_test) # predict the model
207
208 eval_model = r2_score(y_test, y_pred) # this regression evaluate my model to determine if the selected model performance well
209 root_mean_squared_error = np.sqrt(mean_squared_error(y_test, y_pred))
210
211 # Display evaluation metrics
212 print("Model Performance Metrics:")
213 print("-----")
214 print(f"R-squared (R²): {eval_model:.3f}")
215 print(f"Root Mean Squared Error (RMSE): {root_mean_squared_error:.3f}")
216
217 # Perform parameter tuning on the model if needed to improve the performance of your model.
218
219 # Add your parameter tuning code here:
220
221 # Define the base model
222
223 tree = DecisionTreeRegressor(random_state=42)
224
225 # Define the parameter grid
226
227 param_grid = {
228     'max_depth': [3, 5, 7, 9, 12, None],
229     'min_samples_split': [2, 5, 10, 20],
230     'min_samples_leaf': [1, 2, 4, 6],
231     'max_features': [None, 'sqrt', 'log2']
232 }
233
234 # Initialize GridSearchCV
235 grid_search = GridSearchCV(
236     estimator=tree,
237     param_grid=param_grid,
238     cv=5, # 5-fold cross-validation
239     scoring='r2', # Evaluate based on R-squared
240     n_jobs=-1, # Use all CPU cores
241     verbose=1 # Show progress
242 )
243
244 # Fit the grid search on training data
245
246 grid_search.fit(X_train, y_train)
247
248 # Retrieve the best model and parameters
249
250 best_model = grid_search.best_estimator_
251 best_params = grid_search.best_params_
252
253 print("Best Parameters Found:")
254 print(best_params)
255
256 # Evaluate best model on the test data
257
258 y_pred_best = best_model.predict(X_test)
259 r2_best = r2_score(y_test, y_pred_best)
260 rmse_best = np.sqrt(mean_squared_error(y_test, y_pred_best))
261
262 print("\nTuned Model Performance:")
263 print("-----")
264 print(f"R-squared (R²): {r2_best:.3f}")

```

```
265 print(f"Root Mean Squared Error (RMSE): {rmse_best:.3f}")
266 print("-----")
267 # --- Model Prediction ---
268 # TODO: Predict house prices from a new set of feature inputs.
269 # Example new data: Use hypothetical or randomly generated values for the features.
270 #
271 # Example new data: CRIM = 0.2, ZN = 12.5, INDUS = 7.07, CHAS = 0, NOX = 0.5, RM = 6.5, AGE = 68, DIS = 4.0, RAD = 2, TAX =
272
273 # Add your prediction code here:
274
275 # The dictio created below was generated based on the above text with feature names and hypothetical values given.
276
277 new_house_data = {
278     'crime_rate': [0.2],
279     'residential_land_zone': [12.5],
280     'non_retail_business_acres': [7.07],
281     'charles_river_dummy': [0],
282     'nitric_oxides_conc': [0.5],
283     'avg_rooms_per_dwelling': [6.5],
284     'owner_occupied_units_pct': [68],
285     'distance_to_employment': [4.0],
286     'highway_accessibility_index': [2],
287     'property_tax_rate': [250],
288     'pupil_teacher_ratio': [17],
289     'black_population_index': [400],
290     'lower_status_pct': [12]
291 }
292
293 # Convert the dictionary to a DataFrame
294 new_data_df = pd.DataFrame(new_house_data)
295
296 # Use the best model if you tuned, otherwise your base model
297 predicted_value = model.predict(new_data_df) # the best model is being used here to because of the prediction that was made
298
299 print("Predicted Median Home Value:")
300 print(f"${predicted_value[0]:.2f}K")
301
302
```


Inspect the Dataset

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	MEDV	
0	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296.0	15.3	396.90	4.98	24.0	
1	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242.0	17.8	396.90	9.14	21.6	
2	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242.0	17.8	392.83	4.03	34.7	
3	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222.0	18.7	394.63	2.94	33.4	
4	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222.0	18.7	396.90	5.33	36.2	

Dataset Info:

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 506 entries, 0 to 505

Data columns (total 14 columns):

#	Column	Non-Null Count	Dtype
0	CRIM	506 non-null	float64
1	ZN	506 non-null	float64
2	INDUS	506 non-null	float64
3	CHAS	506 non-null	category
4	NOX	506 non-null	float64
5	RM	506 non-null	float64
6	AGE	506 non-null	float64
7	DIS	506 non-null	float64
8	RAD	506 non-null	category
9	TAX	506 non-null	float64
10	PTRATIO	506 non-null	float64
11	B	506 non-null	float64
12	LSTAT	506 non-null	float64
13	MEDV	506 non-null	float64

dtypes: category(2), float64(12)

memory usage: 49.0 KB

Summary Statistics:

	CRIM	ZN	INDUS	NOX	RM	AGE	DIS	TAX	PTRATIO	B	LSTAT
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	3.613524	11.363636	11.136779	0.554695	6.284634	68.574901	3.795043	408.237154	18.455534	356.674032	12.659262
std	8.601545	23.322453	6.860353	0.115878	0.702617	28.148861	2.105710	168.537116	2.164946	91.294864	7.145456
min	0.006320	0.000000	0.460000	0.385000	3.561000	2.900000	1.129600	187.000000	12.600000	0.320000	1.730000
25%	0.082045	0.000000	5.190000	0.449000	5.885500	45.025000	2.100175	279.000000	17.400000	375.377500	6.950000
50%	0.256510	0.000000	9.690000	0.538000	6.208500	77.500000	3.207450	330.000000	19.050000	391.440000	11.360000
75%	3.677083	12.500000	18.100000	0.624000	6.623500	94.075000	5.188425	666.000000	20.200000	396.225000	16.950000
max	88.976200	100.000000	27.740000	0.871000	8.780000	100.000000	12.126500	711.000000	22.000000	396.900000	37.950000

Missing Values:

CRIM	0
ZN	0
INDUS	0
CHAS	0
NOX	0
RM	0
AGE	0
DIS	0
RAD	0
TAX	0
PTRATIO	0
B	0
LSTAT	0
MEDV	0

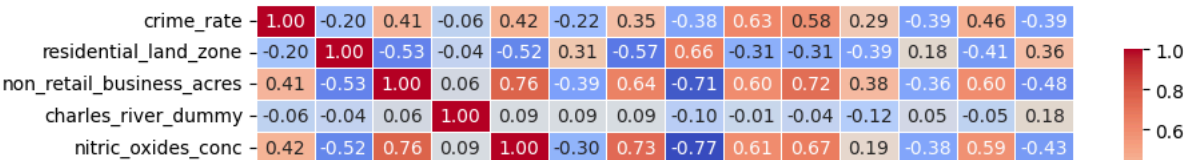
dtype: int64

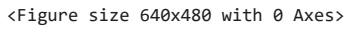
Duplicates Removed

Renamed Columns:

['crime_rate', 'residential_land_zone', 'non_retail_business_acres', 'charles_river_dummy', 'nitric_oxides_conc', 'avg_rooms_per']

Correlation Heatmap





Scatter Plots: Features vs Target (Median Home Value)



Tuned Model Performance:

Selected Features vs Target (Median Home Value)