OPTION #1: Simple Linear Regression in Scikit Learn

Scott Miner

Colorado State University – Global Campus

Abstract

```python
1   import numpy as np
2   import pandas as pd
3   import matplotlib.pyplot as plt
4   import seaborn as sns
5   from sklearn.datasets import load_boston
6   from sklearn.linear_model import LinearRegression, Lasso, Ridge
7   from sklearn.model_selection import train_test_split, cross_val_score
8   from sklearn.metrics import mean_squared_error, make_scorer
9   from sklearn.dummy import DummyRegressor
10  from sklearn.preprocessing import StandardScaler
11  from sklearn.pipeline import make_pipeline
12  from sklearn.utils import shuffle
13  from itertools import combinations
14  #set output options
15
16  params = {'legend.fontsize': 20,
17           'figure.figsize' : (15, 5),
18           'axes.labelsize'  : 30,
19           'axes.titlesize'  : 30,
20           'xtick.labelsize' : 25,
21           'ytick.labelsize' : 25,
22           'lines.markersize': 25,
23           'lines.linewidth' : 5}
24  plt.rcParams.update(params)
25  pd.options.display.width = 0
26  pd.options.display.float_format = '{:,.2f}'.format
27
28  #load data as a bunch object (bo)
29  boston_dataset = load_boston()
30  print(f'Loading boston dataset...')
31  print(f'Type dataset object: {type(boston_dataset)}')
32  print(f'Dataset keys: {boston_dataset.keys()}')
33
34  # create dataframe
35  boston_df = pd.DataFrame(boston_dataset.data, columns=boston_dataset.feature_names)
36  # add target column
37  boston_df['MEDV'] = boston_dataset.target
38  # apply log transformation
39  boston_df['LOGLSTAT'] = boston_df['LSTAT'].apply(np.log)
40
41  # basic dataset informatoin
42  rows, cols = boston_df.shape
43  print(f'Dataset rows: {rows}')
44  print(f'Dataset cols: {cols} ({cols-1} features and 1 target variable)')
45  print(f'Total elements in dataset: {boston_df.size}')
46  print('-------------------------------------------------------------------')
47  # print column descriptions
48  lines = boston_dataset.DESCR.splitlines()[11:28]
49  print('Description of cols:')
50  print(*lines, sep='\n')
51  print('-------------------------------------------------------------------')
52  print('Info:')
53  print(boston_df.info())
54  print('-------------------------------------------------------------------')
55  print(boston_df.describe())
```

*Figure 1.* Program code (part 1)

```python
56  print('-----------------------------------------------------------------')
57  print('First 5 and last 5 rows of the dataset:')
58  print(pd.concat([boston_df.head(), boston_df.tail()]))
59
60  # graphs
61
62  ############### heatmap ###############
63  mask = np.zeros_like(boston_df.corr())
64  mask[np.triu_indices_from(mask)] = True
65
66  ax = sns.heatmap(boston_df.corr().round(2), square=True,
67                   annot=True, mask=mask, cmap='coolwarm', annot_kws={'size': 20})
68  ax.set_title('Heatmap of the Boston Housing Dataset',
69               fontsize=25)
70  ax.tick_params(axis='both', labelsize=20)
71  # color bar object
72  cbar = ax.collections[0].colorbar
73  cbar.ax.tick_params(labelsize=20)
74
75  ############### LSTAT distribution ###############
76  g = sns.displot(boston_df['LSTAT'], bins=30)
77  g.fig.suptitle('Distribution of LSTAT\n(% lower status of the population)',
78                 fontsize=25)
79  g.axes[0,0].set_xlabel('LSTAT', fontsize=20)
80  g.axes[0,0].set_ylabel('Count', fontsize=20)
81  g.axes[0,0].tick_params(axis='both', labelsize=20)
82
83  ############### Histogram: all variables ###############
84  fig, axes = plt.subplots(round(len(boston_df.columns)/3),3,figsize=(20, 40))
85  i = 0
86  for triaxis in axes:
87      for axis in triaxis:
88          if i < len(boston_df.columns):
89              boston_df.hist(column = boston_df.columns[i], bins = 100, ax=axis)
90              axis.set_title(boston_df.columns[i], fontsize=25)
91              axis.tick_params(axis='both', labelsize=20)
92              i += 1
93  fig.suptitle('Histograms for All Variables in the Boston Housing Dataset',
94               fontsize=26)
95
96  ############### Scatterplots ###############
97  features = ['LSTAT', 'RM']
98  target = boston_df['MEDV']
99  fig, axes = plt.subplots(1, len(features), figsize=(10,5))
100
101 for col,ax in zip(features, axes.flat):
102     x = boston_df[col]
103     y = target
104     ax.scatter(x,y,marker='o',s=150)
105     ax.set_title(col, fontsize=26)
106     ax.set_xlabel(col, fontsize=20)
107     ax.set_ylabel('MEDV', fontsize=20)
108     ax.tick_params(axis='both', labelsize=20)
109 fig.suptitle('Scatterplot showing the LSTAT and RM variables against MEDV',
```

No issues found

*Figure 2*. Program code (part 2)

```python
109  fig.suptitle('Scatterplot showing the LSTAT and RM variables against MEDV',
110                   fontsize=30)
111
112  ############### Transformations ###############
113  fig, axes = plt.subplots(1, 2)
114  xs = ['LSTAT', 'LOGLSTAT']
115  colors = ['green', 'red']
116  x_line_points = [[0,40],[0,4]]
117  y_line_points = [[30,0],[50,0]]
118
119  for ax,x,color,x_lp,y_lp in zip(axes.flat,xs,colors,
120                                  x_line_points,y_line_points):
121      ax.scatter(boston_df[x], boston_df['MEDV'], color=color, s=150)
122      ax.set_xlabel(x,fontsize=20)
123      ax.set_ylabel('MEDV',fontsize=20)
124      ax.tick_params(axis='both', labelsize=20)
125      ax.plot(x_lp,y_lp)
126
127  def rms_error(actual, predicted):
128      ' root-mean-squared-error function'
129      # lesser values are better(a<b means a is better)
130      mse = mean_squared_error(actual, predicted)
131      return np.sqrt(mse)
132  rms_scorer = make_scorer(rms_error)
133
134  boston_ftrs = boston_df[['LOGLSTAT', 'RM']]
135  boston_tgt = boston_df['MEDV']
136
137  ############### Train-Test-Split ###############
138  boston_tts = train_test_split(boston_ftrs, boston_tgt,random_state=2021)
139  (boston_train_ftrs, boston_test_ftrs,
140   boston_train_tgt,  boston_test_tgt) = boston_tts
141  print('-----------------------------------------------------------------')
142  print('Split the data into training and test features:' )
143  print(f'Training Features: {boston_train_ftrs.shape}')
144  print(f'Training Target: {boston_train_tgt.shape}')
145  print(f'Testing Features: {boston_test_ftrs.shape}')
146  print(f'Testing Target: {boston_test_tgt.shape}')
147  print('-----------------------------------------------------------------')
148  # preprocessing
149  scaler = StandardScaler()
150  # create models
151  lr = LinearRegression()
152  lasso = Lasso()
153  ridge = Ridge()
154
155  std_lr_pipe = make_pipeline(scaler, lr)
156  std_lasso_pipe = make_pipeline(scaler, lasso)
157  std_ridge_pipe = make_pipeline(scaler, ridge)
158
159  # Create models
160  regressors = {'baseline'       : DummyRegressor(strategy='mean'),
161                'std_lr_pipe'     : std_lr_pipe,
162                'std_lasso_pipe'  : std_lasso_pipe,
163                'std_ridge_pipe'  : std_ridge_pipe}
```

*Figure 3.* Program code (part 3)

```
165  fig, ax = plt.subplots(1, 1, figsize=(8,4))
166  scores = {}
167  for mod_name, model in regressors.items():
168      cv_results = cross_val_score(model,
169                                   boston_train_ftrs, boston_train_tgt,
170                                   scoring = rms_scorer,
171                                   cv=10)
172      key = mod_name
173      scores[key] = [cv_results.mean(), cv_results.std()]
174      lbl = f'{mod_name:s} ({cv_results.mean():5.3f})$\pm${cv_results.std():.2f}'
175      ax.plot(cv_results, 'o--', label=lbl, markersize=11)
176      ax.set_xlabel('CV-Fold #')
177      ax.set_ylabel('RMSE')
178      ax.legend(bbox_to_anchor=(1.00, 1.00), fancybox=True, shadow=True)
179
180  df = pd.DataFrame.from_dict(scores, orient='index').sort_values(0)
181  df.columns = ['RMSE', 'STD_DEV']
182  print('Results on training data')
183  print(df)
184  y_predicted = std_lr_pipe.fit(boston_train_ftrs, boston_train_tgt).predict(boston_test_ftrs)
185  #################### Regression Errors ##############################
186  from itertools import permutations
187  def regression_errors(figsize, predicted, actual, errors='all'):
188      ''' figsize -> subplots;
189          predicted/actual data -> columns in a DataFrame
190          errors -> 'all' or sequence of indices
191      '''
192      fig, axes = plt.subplots(1, 2, figsize=figsize,
193                               sharex=True, sharey=True)
194      df = pd.DataFrame({'actual':actual,
195                         'predicted': predicted})
196
197      for ax, (x,y) in zip(axes, permutations(['actual',
198                                               'predicted'])):
199          # plot the data as '.'; perfect as y=x line
200          ax.plot(df[x], df[y], '.', label='data')
201          ax.plot(df['actual'], df['actual'], '-',
202                  label='perfection')
203          ax.legend()
204
205          ax.set_xlabel('{} Value'.format(x.capitalize()))
206          ax.set_ylabel('{} Value'.format(y.capitalize()))
207          ax.set_aspect('equal')
208
209      axes[1].yaxis.tick_right()
210      axes[1].yaxis.set_label_position('right')
211
212      # show connecting bars from data to perfect
213      # for all or only those specified?
214      if errors == 'all':
215          errors = range(len(df))
216      if errors:
217          acts = df.actual.iloc[errors]
```

Figure 4. Program code (part 4)

```
217          acts = df.actual.iloc[errors]
218          preds = df.predicted.iloc[errors]
219          axes[0].vlines(acts, preds, acts, 'r')
220          axes[1].hlines(acts, preds, acts, 'r')
221
222  regression_errors((10, 5), y_predicted, boston_test_tgt, errors=[0,1,2,3,4,
223                                                       122,123,124,125,126])
224
225  def regression_residuals(ax, predicted, actual,
226                          show_errors=None, right=False):
227      ''' figsize -> subplots;
228          predicted/actual data -> columns of a DataFrame
229          errors -> 'all' or sequence of indices
230      '''
231      df = pd.DataFrame({'actual':actual,
232                          'predicted':predicted})
233      df['error'] = df.actual - df.predicted
234      ax.plot(df.predicted, df.error, '.')
235      ax.plot(df.predicted, np.zeros_like(predicted), '-')
236
237      if right:
238          ax.yaxis.tick_right()
239          ax.yaxis.set_label_position('right')
240          ax.set_xlabel('Predicted Value')
241          ax.set_ylabel('Residual')
242
243      if show_errors == 'all':
244          show_errors = range(len(df))
245      if show_errors:
246          preds = df.predicted.iloc[show_errors]
247          errors = df.error.iloc[show_errors]
248          ax.vlines(preds, 0, errors, 'r')
249
250  fig, ax= plt.subplots(1,1,figsize=(8,8))
251
252  regression_residuals(ax, y_predicted, boston_test_tgt, show_errors=[0,1,2,3,4,
253                                                       122,123,124,125,126])
254  ax.set_xlabel('Predicted')
255  ax.set_ylabel('Residual')
256
257
258  coeff = list(zip(boston_ftrs, std_lr_pipe.named_steps['linearregression'].coef_))
259
260  df_coeff = pd.DataFrame(coeff)
261  df_coeff.columns = ['variable', 'coeff']
262  df_coeff.set_index('variable')
263  df_intercept = pd.DataFrame([['INTERCEPT', std_lr_pipe.named_steps['linearregression'].intercept_]])
264  df_intercept.columns = ['variable', 'coeff']
265  df_results = df_intercept.append(df_coeff, ignore_index=True)
266  df_results = df_results.set_index(['variable'])
267  df_results.index.name=None
268  df_results
269  print('----------------------------------------------------------------------')
270  print('RMSE on hold-out data:')
271  test_score = np.sqrt(mean_squared_error(boston_test_tgt, y_predicted))
```

*Figure 5.* Program code (part 5)

```
269  print('------------------------------------------------------------------------')
270  print('RMSE on hold-out data:')
271  test_score = np.sqrt(mean_squared_error(boston_test_tgt, y_predicted))
272  test_results = pd.DataFrame([('std_lr_pipe', test_score)])
273  test_results.columns = ['pipeline', 'RMSE']
274  test_results = test_results.set_index('pipeline')
275  test_results.index.name=None
276  print(test_results)
277  print('------------------------------------------------------------------------')
278  print('Coefficients:')
279  print(df_results)
280  print('------------------------------------------------------------------------')
281  print('Display True vs. Actual Values:')
282  y_test = pd.DataFrame(boston_test_tgt)
283  y_test['PREDS'] = y_predicted
284  y_test = y_test.rename(columns={'MEDV': 'TRUE'})
285  y_test = y_test[['PREDS', 'TRUE']]
286  print(y_test)
287
288  plt.show()
```

*Figure 6.* Program code (part 6)



*Figure 7.* Program output (part 1)

```
------------------------------------------------------------
Split the data into training and test features:
Training Features: (379, 2)
Training Target: (379,)
Testing Features: (127, 2)
Testing Target: (127,)
------------------------------------------------------------
Results on training data
          RMSE   STD_DEV
ridge     4.93     0.86
lr        4.93     0.85
lasso     5.20     1.02
baseline  9.47     1.52
------------------------------------------------------------
RMSE on hold-out data:
    RMSE
lr  5.05
------------------------------------------------------------
Coefficients:
           coeff
INTERCEPT  24.52
LOGLSTAT  -10.35
RM          3.58
------------------------------------------------------------
Display True vs. Actual Values:
      LSTAT    RM   PREDS   TRUE
210   17.27 5.96   16.37 21.70
24    16.30 5.92   16.84 15.60
36    11.41 5.84   20.24 20.00
439   22.88 5.63   12.27 12.80
161    1.73 7.49   45.66 50.00
..      ...   ...     ...    ...
336    9.80 5.87   21.91 19.50
385   30.81 5.28    7.94  7.20
352    7.79 5.88   24.34 18.60
106   18.66 5.84   15.13 19.50
391   18.76 6.05   15.84 23.20
```
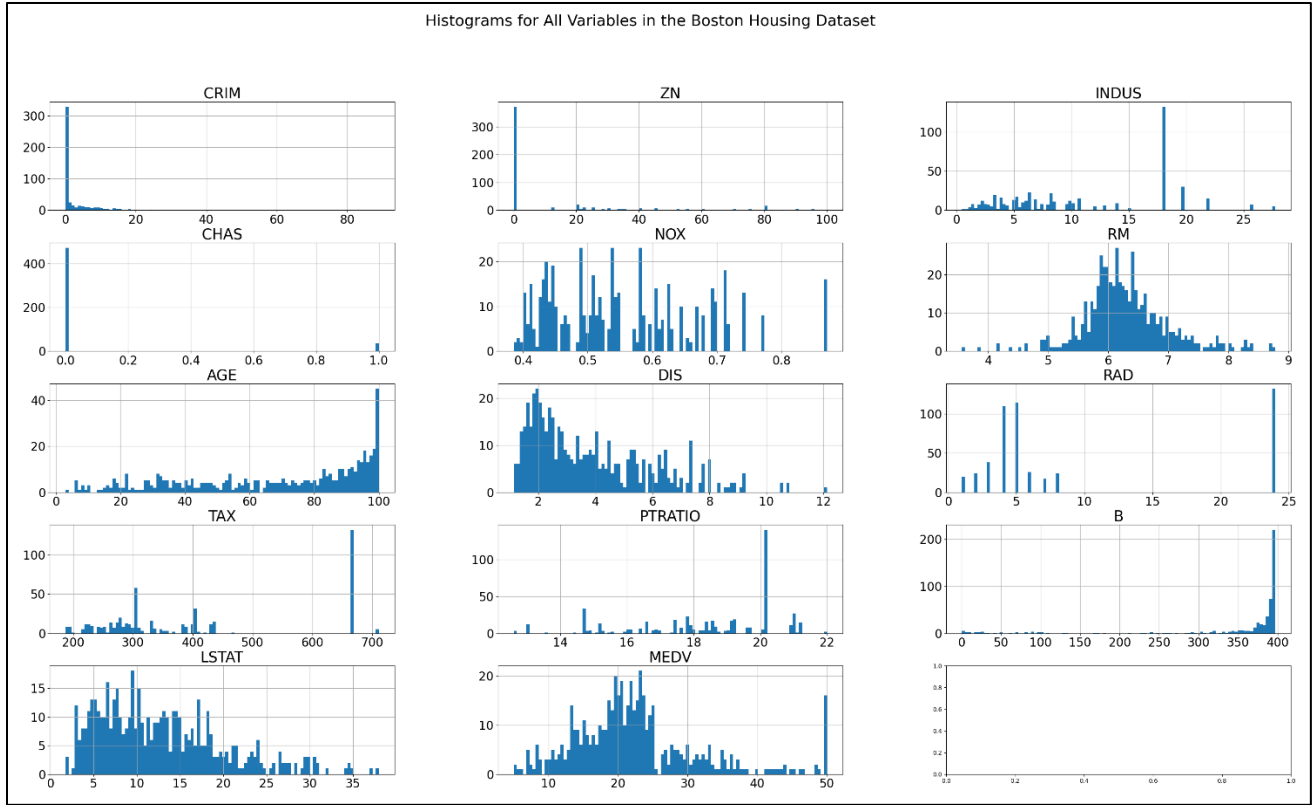
*Figure 8.* Program output (part 2)

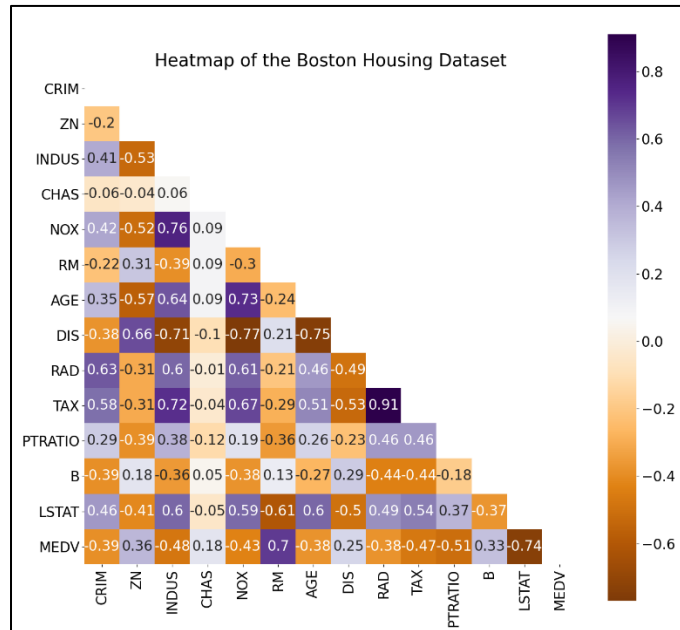*Figure 9*. Histogram of variables in the Boston housing dataset
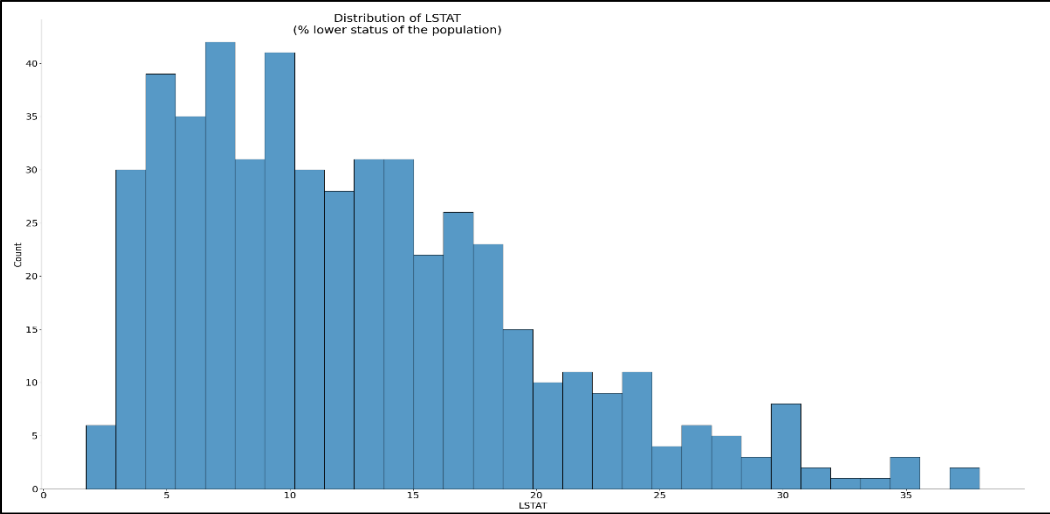


*Figure 10*. Heatmap of the Boston housing dataset
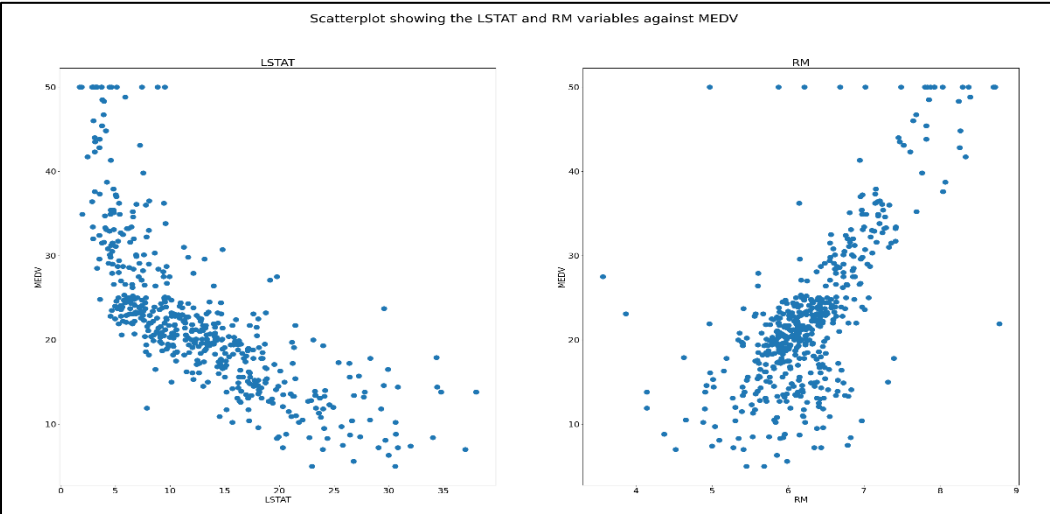
*Figure 11.* Distribution of the *LSTAT* variable



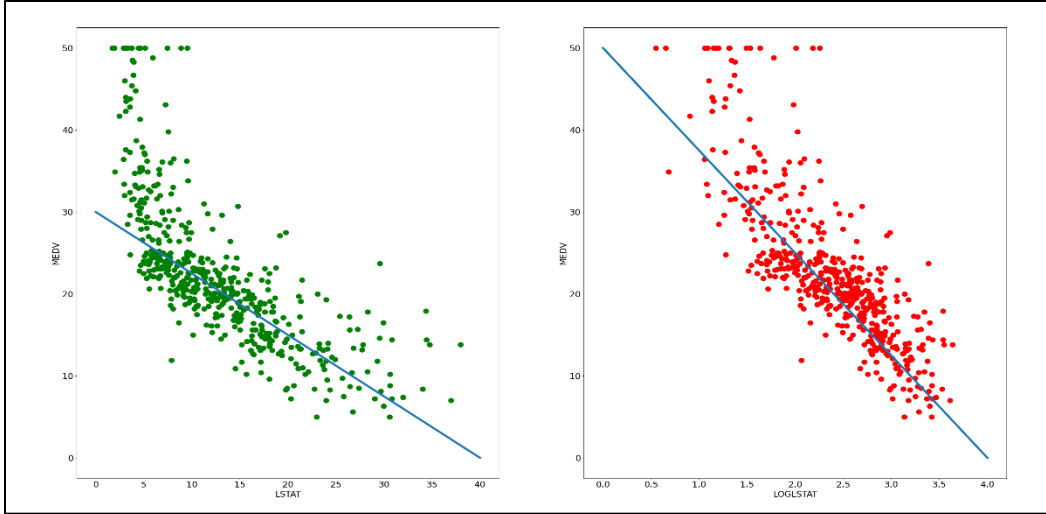*Figure 12.* Scatterplots showing the *LSTAT* and *RM* variables against *MEDV*

*Figure 13.* Performing a logarithmic transformation of the *LSTAT* variable
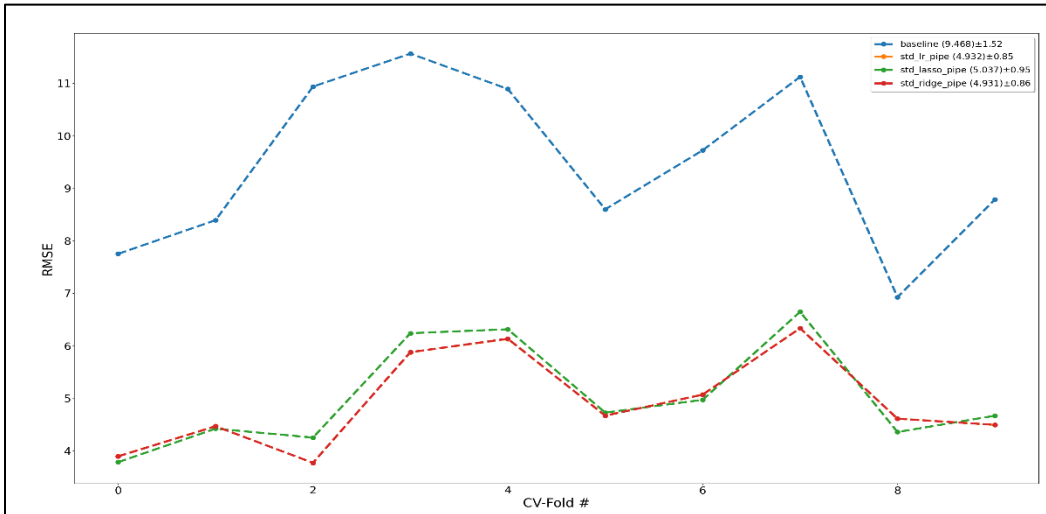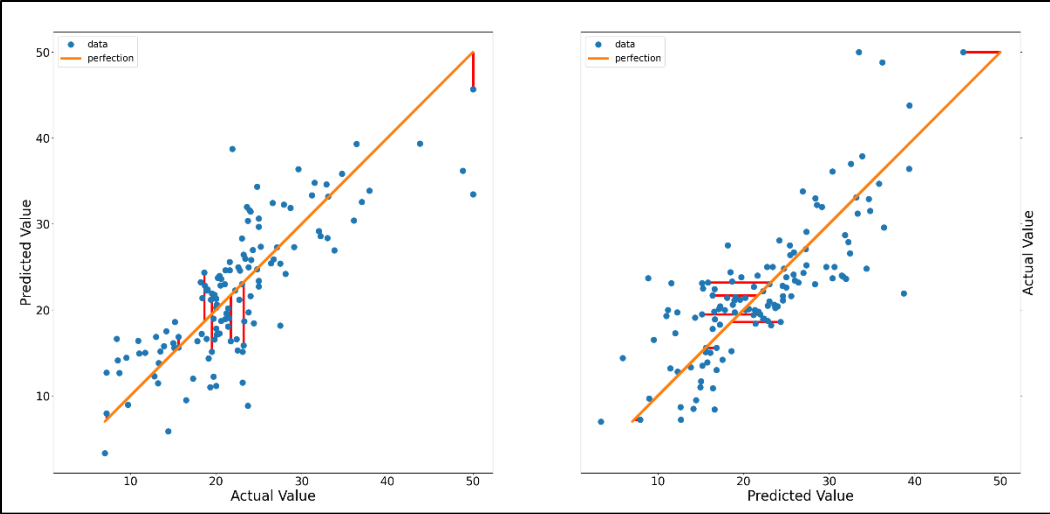


*Figure 14.* Model evaluation on 10-fold cross-validation
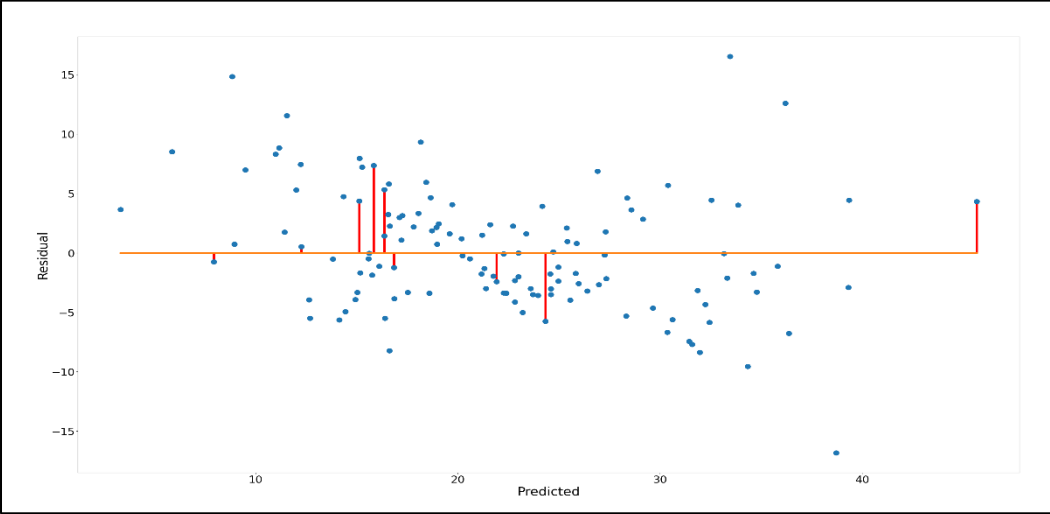
*Figure 15.* Regression Errors



*Figure 16.* Residual Plot

**OPTION #1: Simple Linear Regression in Scikit Learn**

This paper describes a linear regression model in Python that predicts the median value of an owner-occupied home in the $1,000s from the Boston housing dataset. The dataset contains 506 rows and 14 features, totaling 7,804 values. Figure 7, the first half of the program's console output, shows each dataset feature's name and the corresponding description. Figures 1 – 6 show the code to create the program and its output. There are no null values in the dataset, and all the variables are floating-point numbers that occupy 64 bits in memory.

**Exploratory Data Analysis (EDA)**

The program outputs descriptive statistics for all the variables and the first and last five rows of the dataset. Figure 9 shows a histogram for all dataset variables. Agarwal (2018) writes that the target variable, *MEDV*, which indicates the median value of an owner-occupied home in the $1,000s, appears normally distributed with a few outliers. Figure 10 shows a correlation matrix and any multicollinearity that exists between variables. When performing linear regression, we seek to identify variables that strongly correlate with the target variable. The correlation matrix shows that *RM*, the average number of rooms per house, correlates positively with the target variable ($r=0.7$) and that *LSTAT*, the percent of the lower status of the population, correlates negatively with the target ($r=-0.74$). Therefore, our linear regression model will focus on these two variables and ignore the remainder since they do not strongly correlate with the target variable.

Figure 11 takes a closer look at the distribution of the *LSTAT* variable. In contrast to the distribution of the *RM* variable, which is normally distributed, the distribution of the *LSTAT* variable appears positively skewed—it has a tail on the right side. Figure 12 plots the *RM* and *LSTAT* variables against the *MEDV* variable. The data in the scatterplot on the left-hand side of the image, showing *LSTAT* against *MEDV,* appears to follow more of a curved distribution than

the plot of *RM* against *MEDV*, shown on the figure's right-hand side. Therefore, Chung (2019) writes that we can transform the *LSTAT* variable logarithmically to prevent the model from underfitting.

Underfitting occurs when the model is incapable of capturing the variability of the training data (Jabbar & Khan, 2014). Chung states that by logarithmically transforming the *LSTAT* variable, we minimize the nonlinear relationship in the data and create a more accurate model. Figure 13 compares the *LSTAT* variable with its logarithmically transformed counterpart, *LOGLSTAT*, plotted against *MEDV*. We capture more data by drawing a diagonal line through the right-hand image, which contains the logarithmically transformed variable, rather than the left-hand image containing the untransformed data. Therefore, using this transformed variable to train our linear regression model allows it to capture more of the training data's variability.

**Model Training, Evaluation, and Results**

The program divides the training features, which contain just the *LOGLSTAT* and *RM* variables, and the target variable, *MEDV*, into 75% and 25% training and testing sets, respectively. The training set contains 379 rows, and the testing set contains 127 rows, as shown by the program's console output in Figure 8. The program then computes the root mean squared error (RMSE) for a baseline regressor that uses the target variable's mean as its prediction and three linear regression models, (a) good old-fashioned (GOF) linear regression, (b) lasso regression, and (c) ridge regression, using 10-fold cross-validation over the training data.

Figure 14 shows the RMSE for each of the models across each of the ten folds. The lower the RMSE, the better the model. The image shows that GOF linear regression ($RMSE =$ $4.932 \pm 0.85$) and ridge regression ($RMSE = 4.931 \pm 0.86$) performed nearly identically. Therefore, GOF linear regression was chosen to make predictions on the test data. The model achieved an RMSE of 5.05 on the test data, as shown in Figure 8, using the following equation:

$\widehat{Y}_i = -10.35 \times \ln(LSTAT_i) + 3.58 \times RM_i + 24.52$. The program outputs the independent

variables and their corresponding coefficients to the console, as shown in Figure 8.

**Discussion and Predictions**

How can we interpret this output? The RMSE values of 4.93 and 5.05, which the GOF

linear regression model achieved on the training and testing sets, respectively, can be interpreted

as saying that the model is roughly $5,000 off the actual value, on average. Additionally, in the

equation that generates the model's predictions, we can interpret the coefficient of the *LSTAT*

variable, the percentage of the lower status of the population, as saying that a 1% increase in

*LSTAT decreases* the median value of an owner-occupied home by about $\frac{10.35}{100} = .1035$ or

$103.5, holding all other independent variables constant (Chung, 2019; *Interpreting Log*

*Transformations in a Linear Model | University of Virginia Library Research Data Services +*

*Sciences*, n.d.). Further, Chung (2019) writes that we can interpret the coefficient of the *RM*

variable, the average number of rooms per home, as saying that for every one-unit increase in a

house's average number of rooms, the median value increases by about $3,580, holding all other

independent variables constant. Lastly, the y-intercept indicates that the starting price of a house

in Boston in 1979 would be around $24,520.

The program's final output, shown in Figure 8, displays ten records from the test data set

and their corresponding *LSTAT, RM,* predicted, and true values. One can plug the values from

the *LSTAT* and *RM* variables into the equation above to understand how the model generates its

predictions. Furthermore, Figures 15 and 16, which plot the *actual vs. predicted values* and the

*predicted vs. residual values* of the test dataset, respectively, display these same ten records with

their predictive errors highlighted in red. For instance, residual plots indicate what we need to do

to fix our predictions (Fenner, 2019). From Figure 16, we can see that the model consistently

under-predicts a few records around the $15,000 mark by about $5,000, as can be verified with the program's console output, shown in Figure 8.

## Conclusion

In conclusion, this paper described a simple linear regression model in Python, developed using scikit-learn, to predict the median value of an owner-occupied home in the $1,000s from the Boston housing dataset. The paper gave a brief overview of the dataset and used graphical approaches, including heatmaps, to explore the data to determine variables strongly correlated with the target variable. The paper also discussed techniques to logarithmically transform the independent variables to avoid underfitting the model and used graphical output to verify the correctness of these techniques. The best performing linear regression model was selected to create predictions on the test data and achieved an RMSE of 5.05. Finally, the paper discussed the implications of this model and its predictions and used graphical output to support these discussions.

References

Agarwal, A. (2018, October 5). *Linear Regression on Boston Housing Dataset*. Medium.

> https://towardsdatascience.com/linear-regression-on-boston-housing-dataset-
>
> f409b7e4a155

Chung, C. (2019, August 29). *Used Linear Regression to Model and Predict Housing Prices with*

> *The Classic Boston Housing Dataset*. Chris Chung.
>
> https://chrispfchung.github.io//model%20data/boston-housing-data/

Fenner, M. (2019). *Machine Learning in Python for Everyone*. Addison-Wesley.

*Interpreting Log Transformations in a Linear Model | University of Virginia Library Research*

> *Data Services + Sciences*. (n.d.). Retrieved September 5, 2021, from
>
> https://data.library.virginia.edu/interpreting-log-transformations-in-a-linear-model/

Jabbar, H. K., & Khan, R. Z. (2014). Methods to Avoid Over-Fitting and Under-Fitting in

> Supervised Machine Learning (Comparative Study). *Computer Science, Communication,*
>
> *and Instrumentation Devices*, 163–172. https://doi.org/10.3850/978-981-09-5247-1_017