

## Building a Linear Regression Model to Predict Diabetes Progression

```
import sys # Read system parameters
import numpy as np # Work with multi-dimensional arrays
import pandas as pd # Manipulate and analyze data
```

```

import matplotlib.cm as cm
import seaborn as sns
import sklearn
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

%matplotlib inline

# Summarize software libraries used.
print('Libraries used in this project:')
print('- Python ({}).format(sys.version)')
print('- Numpy ({}).format(np.__version__)')
print('- pandas ({}).format(pd.__version__)')
print('- Matplotlib ({}).format(mpl.__version__)')
print('- Seaborn ({}).format(sns.__version__)')
print('- scikit-learn ({}).format(sklearn.__version__)')

# Load the dataset.
diabetes = datasets.load_diabetes()
print('Loaded {} records.'.format(len(diabetes.data)))

Libraries used in this project:
- Python 3.8.5 (default, Sep 3 2020, 21:29:08) [MSC v.1916 64 bit (AMD64)]
- Numpy 1.19.2
- pandas 1.1.3
- Matplotlib 3.3.2
- Seaborn 0.11.0
- scikit-learn 0.23.2

Loaded 442 records.

Get acquainted with the dataset

```

---

```

In [2]: # Convert array to pandas DataFrame.

# View data types and see if there are missing entries.

# View first 10 records.

In [3]: print(diabetes)

'data': array([[ 0.03807591,  0.05068012,  0.06169621, ..., -0.00259226,
  0.0190842, -0.01764613],
 [-0.00188202, -0.04646164, -0.05147406, ..., -0.03949338,
  -0.06832374, -0.09220405],
 [ 0.08529891,  0.05068012,  0.04445121, ..., -0.00259226,
  0.00286377, -0.02593034],
 ...,
 [ 0.04170844,  0.05068012, -0.01590626, ..., -0.01107952,
  -0.04687948,  0.01549073],
 [-0.04547248, -0.04646164,  0.03906215, ...,  0.02655962,
  0.04452837, -0.02593034],
 [-0.04547248, -0.04646164, -0.0730303, ..., -0.03949338,
  -0.00421986,  0.00306441]), 'target': array([151., 75., 141., 206., 135., 97., 138., 63., 110.,
  310., 101.,
  69., 179., 185., 118., 171., 166., 144., 97., 168., 65., 49.,
  68., 245., 184., 202., 137., 85., 131., 283., 129., 59., 341.,
  87., 65., 102., 265., 276., 252., 90., 100., 55., 61., 92.,
  258., 53., 100., 192., 75., 102., 155., 205., 58., 104., 163.]

```

128.	52.	37.	170.	170.	61.	144.	52.
150.	97.	160.	178.	48.	270.	202.	111.
200.	252.	113.	143.	51.	52.	210.	65.
42.	111.	98.	164.	48.	96.	90.	162.
83.	128.	102.	302.	198.	95.	53.	134.

```

173, 180, 84, 121, 161, 95, 109, 115, 268, 274, 158,
197, 83, 103, 275, 85, 280, 336, 281, 274, 235,
60, 174, 259, 178, 128, 96, 126, 288, 88, 292, 71,
197, 186, 25, 84, 96, 195, 53, 217, 172, 131, 214,
59, 70, 220, 268, 152, 47, 74, 295, 101, 151, 127,
237, 225, 81, 151, 107, 64, 138, 185, 265, 101, 137,
143, 141, 79, 292, 178, 91, 116, 86, 122, 72, 129,
142, 90, 158, 39, 196, 222, 277, 99, 196, 202, 155,
77, 191, 70, 73, 49, 65, 263, 248, 96, 218, 185,
78, 93, 252, 150, 77, 208, 77, 108, 160, 53, 220,
154, 259, 90, 246, 124, 67, 72, 257, 262, 275, 177,
71, 47, 197, 78, 98, 208, 215, 303, 245, 91,
150, 310, 153, 346, 63, 89, 50, 39, 103, 308, 116,
145, 74, 45, 115, 264, 87, 202, 127, 182, 241, 66,
94, 283, 64, 102, 200, 265, 94, 230, 181, 156, 233,
60, 219, 80, 69, 332, 248, 84, 200, 55, 85, 89,
31, 129, 83, 275, 65, 198, 236, 253, 124, 44, 172,
114, 142, 109, 180, 144, 163, 147, 97, 220, 190, 109,
91, 122, 230, 242, 248, 249, 192, 131, 237, 78, 135,
244, 199, 270, 164, 72, 96, 306, 91, 214, 95, 216,
263, 178, 113, 200, 139, 139, 88, 148, 88, 243, 71,
77, 109, 272, 60, 54, 221, 90, 311, 281, 182, 321,
58, 262, 206, 233, 242, 123, 167, 63, 197, 71, 168,
140, 217, 121, 235, 245, 40, 52, 104, 132, 88, 69,
218, 72, 201, 110, 51, 277, 63, 118, 69, 273, 235,
43, 198, 242, 232, 175, 93, 168, 275, 296, 218, 175,
140, 189, 181, 209, 136, 261, 113, 131, 174, 257, 55,
84, 42, 146, 212, 233, 91, 111, 152, 120, 67, 310,
94, 183, 66, 173, 72, 49, 64, 48, 178, 104, 132,
220, 57, j), 'frame': None, 'DESCR': '. diabetes dataset\n\n-----\n\n
N'ten baseline variables, age, sex, mass index, average blood pressure, and six blood sugar measures
to be obtained for each of n=442 diabetes patients, as well as the response of interest. Quantitative
easure of disease progression one year after baseline.\n\nData Set Characteristics:\n\n
Number of instances: 442\n
Number of Attributes: First 10 columns are numeric predictive values\n
Target: Column n-1 is a quantitative measure of disease progression one year after baseline\n\n
Attribute Information:\n
1 age in years\n
2 sex\n
3 bmi, high-density lipoprotein\n
4 tc, T-cells (a type of white blood cells)\n
5 ldl, low-density lipoprotein\n
6 avg, average blood pressure\n
7 ln10, log10-transformed\n
8 ln10, log10-transformed\n
9 ln10, log10-transformed\n
10 ln10, log10-transformed\n
11 ln10, log10-transformed\n
12 ln10, log10-transformed\n
13 ln10, log10-transformed\n
14 ln10, log10-transformed\n
15 ln10, log10-transformed\n
16 ln10, log10-transformed\n
17 ln10, log10-transformed\n
18 ln10, log10-transformed\n
19 ln10, log10-transformed\n
20 ln10, log10-transformed\n
21 ln10, log10-transformed\n
22 ln10, log10-transformed\n
23 ln10, log10-transformed\n
24 ln10, log10-transformed\n
25 ln10, log10-transformed\n
26 ln10, log10-transformed\n
27 ln10, log10-transformed\n
28 ln10, log10-transformed\n
29 ln10, log10-transformed\n
30 ln10, log10-transformed\n
31 ln10, log10-transformed\n
32 ln10, log10-transformed\n
33 ln10, log10-transformed\n
34 ln10, log10-transformed\n
35 ln10, log10-transformed\n
36 ln10, log10-transformed\n
37 ln10, log10-transformed\n
38 ln10, log10-transformed\n
39 ln10, log10-transformed\n
40 ln10, log10-transformed\n
41 ln10, log10-transformed\n
42 ln10, log10-transformed\n
43 ln10, log10-transformed\n
44 ln10, log10-transformed\n
45 ln10, log10-transformed\n
46 ln10, log10-transformed\n
47 ln10, log10-transformed\n
48 ln10, log10-transformed\n
49 ln10, log10-transformed\n
50 ln10, log10-transformed\n
51 ln10, log10-transformed\n
52 ln10, log10-transformed\n
53 ln10, log10-transformed\n
54 ln10, log10-transformed\n
55 ln10, log10-transformed\n
56 ln10, log10-transformed\n
57 ln10, log10-transformed\n
58 ln10, log10-transformed\n
59 ln10, log10-transformed\n
60 ln10, log10-transformed\n
61 ln10, log10-transformed\n
62 ln10, log10-transformed\n
63 ln10, log10-transformed\n
64 ln10, log10-transformed\n
65 ln10, log10-transformed\n
66 ln10, log10-transformed\n
67 ln10, log10-transformed\n
68 ln10, log10-transformed\n
69 ln10, log10-transformed\n
70 ln10, log10-transformed\n
71 ln10, log10-transformed\n
72 ln10, log10-transformed\n
73 ln10, log10-transformed\n
74 ln10, log10-transformed\n
75 ln10, log10-transformed\n
76 ln10, log10-transformed\n
77 ln10, log10-transformed\n
78 ln10, log10-transformed\n
79 ln10, log10-transformed\n
80 ln10, log10-transformed\n
81 ln10, log10-transformed\n
82 ln10, log10-transformed\n
83 ln10, log10-transformed\n
84 ln10, log10-transformed\n
85 ln10, log10-transformed\n
86 ln10, log10-transformed\n
87 ln10, log10-transformed\n
88 ln10, log10-transformed\n
89 ln10, log10-transformed\n
90 ln10, log10-transformed\n
91 ln10, log10-transformed\n
92 ln10, log10-transformed\n
93 ln10, log10-transformed\n
94 ln10, log10-transformed\n
95 ln10, log10-transformed\n
96 ln10, log10-transformed\n
97 ln10, log10-transformed\n
98 ln10, log10-transformed\n
99 ln10, log10-transformed\n
100 ln10, log10-transformed\n
101 ln10, log10-transformed\n
102 ln10, log10-transformed\n
103 ln10, log10-transformed\n
104 ln10, log10-transformed\n
105 ln10, log10-transformed\n
106 ln10, log10-transformed\n
107 ln10, log10-transformed\n
108 ln10, log10-transformed\n
109 ln10, log10-transformed\n
110 ln10, log10-transformed\n
111 ln10, log10-transformed\n
112 ln10, log10-transformed\n
113 ln10, log10-transformed\n
114 ln10, log10-transformed\n
115 ln10, log10-transformed\n
116 ln10, log10-transformed\n
117 ln10, log10-transformed\n
118 ln10, log10-transformed\n
119 ln10, log10-transformed\n
120 ln10, log10-transformed\n
121 ln10, log10-transformed\n
122 ln10, log10-transformed\n
123 ln10, log10-transformed\n
124 ln10, log10-transformed\n
125 ln10, log10-transformed\n
126 ln10, log10-transformed\n
127 ln10, log10-transformed\n
128 ln10, log10-transformed\n
129 ln10, log10-transformed\n
130 ln10, log10-transformed\n
131 ln10, log10-transformed\n
132 ln10, log10-transformed\n
133 ln10, log10-transformed\n
134 ln10, log10-transformed\n
135 ln10, log10-transformed\n
136 ln10, log10-transformed\n
137 ln10, log10-transformed\n
138 ln10, log10-transformed\n
139 ln10, log10-transformed\n
140 ln10, log10-transformed\n
141 ln10, log10-transformed\n
142 ln10, log10-transformed\n
143 ln10, log10-transformed\n
144 ln10, log10-transformed\n
145 ln10, log10-transformed\n
146 ln10, log10-transformed\n
147 ln10, log10-transformed\n
148 ln10, log10-transformed\n
149 ln10, log10-transformed\n
150 ln10, log10-transformed\n
151 ln10, log10-transformed\n
152 ln10, log10-transformed\n
153 ln10, log10-transformed\n
154 ln10, log10-transformed\n
155 ln10, log10-transformed\n
156 ln10, log10-transformed\n
157 ln10, log10-transformed\n
158 ln10, log10-transformed\n
159 ln10, log10-transformed\n
160 ln10, log10-transformed\n
161 ln10, log10-transformed\n
162 ln10, log10-transformed\n
163 ln10, log10-transformed\n
164 ln10, log10-transformed\n
165 ln10, log10-transformed\n
166 ln10, log10-transformed\n
167 ln10, log10-transformed\n
168 ln10, log10-transformed\n
169 ln10, log10-transformed\n
170 ln10, log10-transformed\n
171 ln10, log10-transformed\n
172 ln10, log10-transformed\n
173 ln10, log10-transformed\n
174 ln10, log10-transformed\n
175 ln10, log10-transformed\n
176 ln10, log10-transformed\n
177 ln10, log10-transformed\n
178 ln10, log10-transformed\n
179 ln10, log10-transformed\n
180 ln10, log10-transformed\n
181 ln10, log10-transformed\n
182 ln10, log10-transformed\n
183 ln10, log10-transformed\n
18
```

```

1  -0.001882  -0.044642  -0.051474  -0.026328  -0.008449  -0.019163  0.074412  -0.039493  -0.068340  -0.092204  75.0
2  0.085299  0.050680  0.044451  -0.005671  -0.045599  -0.034194  -0.032356  -0.002592  0.002864  -0.025930  141.0
3  -0.089063  -0.044642  -0.011595  -0.036656  0.012191  0.024991  -0.036038  0.034309  0.022692  -0.009362  206.0
4  0.005383  -0.044642  -0.036385  0.021872  0.003935  0.015596  0.008142  -0.002592  -0.031991  -0.046641  135.0
5  -0.092695  -0.044642  -0.040696  -0.019442  -0.068991  -0.079288  0.041277  -0.076395  -0.041180  -0.096346  97.0
6  -0.045472  0.050680  -0.047163  -0.015999  -0.040096  -0.024800  0.000779  -0.039493  -0.062913  -0.038357  138.0
7  0.063504  0.050680  -0.001895  0.066630  0.090620  0.108914  0.022869  0.017703  -0.035817  0.003064  63.0
8  0.041708  0.050680  0.061696  -0.040099  -0.013953  0.006202  -0.028674  -0.002592  -0.014956  0.011349  110.0
9  -0.070900  -0.044642  0.039062  -0.033214  -0.012577  -0.034508  -0.024993  -0.002592  0.067736  0.013504  310.0

```

## Examine the distribution of various features

```

In [11]: # Use Matplotlib to plot distribution histograms for all features.

import matplotlib.pyplot as plt

plt.rcParams['axes', labelsize=14]
plt.rcParams['xtick', labelsize=12]
plt.rcParams['ytick', labelsize=12]

df.hist(figsize=(20,20))
plt.show()

```

The figure displays 12 histograms arranged in a 4x3 grid, showing the distribution of various features. The features are labeled as follows:

- age**: Distribution of age values, ranging from approximately -10 to 10.
- sex**: Distribution of sex values, with a peak at -0.04 and another at 0.04.
- bmi**: Distribution of BMI values, ranging from approximately -10 to 15.
- tpo**: Distribution of tpo values, ranging from approximately -10 to 10.
- s1**: Distribution of s1 values, ranging from approximately -10 to 15.
- s2**: Distribution of s2 values, ranging from approximately -10 to 20.
- s3**: Distribution of s3 values, ranging from approximately -10 to 15.
- s4**: Distribution of s4 values, ranging from approximately -5 to 15.
- s5**: Distribution of s5 values, ranging from approximately -10 to 10.
- s6**: Distribution of s6 values, ranging from approximately -15 to 10.
- target**: Distribution of target values, ranging from approximately 50 to 350.
- unlabeled**: Distribution of an unlabeled feature, ranging from approximately -15 to 10.

## Examine a general summary of statistics

```

In [12]: # View summary statistics (mean, standard deviation, min, max, etc.) for each feature.
df.describe()

```

	age	sex	bmi	bp	s1	s2
count	4.420000e+02	4.420000e+02	4.420000e+02	4.420000e+02	4.420000e+02	4.420000e+02

```

std      4.761905e-02      4.761905e-02      4.761905e-02      4.761905e-02      4.761905e-02      4.761905e-02      4.761905e-02      4.761905e-02
min      -1.072256e-01      -4.464164e-02      -9.072530e-03      -1.123996e-02      -1.267807e-02      -1.156131e-02      -1.023071e-02      -7.639450e-02      -1.260974e-01
      01              01              01              01              01              01              01              02              01
25%      -3.729927e-02      -4.464164e-02      -3.422907e-02      -3.665645e-02      -3.424748e-02      -3.0335840e-02      -3.511716e-02      -3.949338e-02      -3.324879e-02
      02              02              02              02              02              02              02              02              02
50%      5.383060e-03      -4.464164e-02      -7.283766e-02      5.670611e-03      -4.320866e-02      -3.819056e-02      -6.584468e-02      -2.592262e-02      -1.947634e-03
      03              03              03              03              03              03              03              03              03
75%      3.807591e-02      5.068012e-02      3.124802e-02      3.564384e-02      2.835801e-02      2.984439e-02      2.931150e-02      3.408866e-02      3.243323e-02
      04              04              04              04              04              04              04              04              04
max      1.107267e-01      5.068012e-02      1.705552e-01      1.320442e-01      1.539137e-01      1.887806e-01      1.811791e-01      1.852344e-01      1.335990e-01
      05              06              07              08              09              10              11              12              13

# View the correlation values for each feature compared to the label.

df.corr()

age      sex      bmi      bp      s1      s2      s3      s4      s5      s6      target
age      1.000000      0.173737      0.185085      0.335427      0.260601      0.219243      -0.075181      0.203841      0.270777      0.301731      0.187889
sex      0.173737      1.000000      0.088161      0.241013      0.035277      0.142637      -0.379090      0.032115      0.149918      0.028133      0.043062
bmi      0.185085      0.088161      1.000000      0.395415      0.249777      0.261170      -0.366811      0.413807      0.446159      0.388680      0.586450
bp      0.335427      0.241013      0.395415      1.000000      0.242470      0.185558      -0.178761      0.257653      0.393478      0.390429      0.441484
s1      0.260601      0.035277      0.249777      0.242470      1.000000      0.896663      0.051519      0.542207      0.515501      0.325717      0.212022
s2      0.219243      0.142637      0.261170      0.185558      0.896663      1.000000      -0.196455      0.659817      0.318353      0.290600      0.174054
s3      -0.075181      -0.379090      -0.366811      -0.178761      0.051519      -0.196455      1.000000      -0.738493      -0.398577      -0.273697      -0.394789
s4      0.203841      0.032115      0.413807      0.257653      0.542207      0.659817      -0.738493      1.000000      0.617857      0.471712      0.430453
s5      0.270777      0.149918      0.446159      0.393478      0.515501      0.318353      -0.398577      0.617857      1.000000      0.464670      0.565883
s6      0.301731      0.028133      0.388680      0.390429      0.325717      0.290600      -0.273697      0.471712      0.464670      1.000000      0.382483
target   0.187889      0.043062      0.586450      0.441484      0.212022      0.174054      -0.394789      0.430453      0.565883      0.382483      1.000000

df.corr()["target"].sort_values()

s3      -0.394789
sex      0.043062
s2      0.174054
age      0.187889
s1      0.212022
s6      0.382483
s4      0.430453
bp      0.441484
s5      0.565883
bmi      0.586450
target   1.000000
Name: target, dtype: float64

Split the label from the dataset

df = pd.read_csv("diabetes.csv")

# Split the training and test datasets and their labels.

# Compare the number of rows and columns in the original data to the training and test sets.

df.shape

(442, 11)

X = df.iloc[:,0:10]
y = df.iloc[:,10]

X.values, y.values

(array([[ 0.3807591,  0.05068012,  0.06169621, ..., -0.00259226,
         0.01990842, -0.01764613],
        [-0.00185202, -0.0466164, -0.05147406, ..., -0.03949338,
         -0.06832974, -0.09224005],
        [ 0.08529891,  0.05068012,  0.04445121, ..., -0.00259226,
         0.00286377, -0.02593034],
        ...,
        [ 0.04170844,  0.05068012, -0.01590626, ..., -0.01070952,
         -0.04687948,  0.01549073],
        [-0.04547248, -0.0466164,  0.0396215, ...,  0.02655962,
         0.04452837, -0.02593034],
        [-0.04547248, -0.0466164, -0.0730303, ..., -0.03949338,
         -0.00421986,  0.03036411]],
        array([ 75, 141, 206, 135, 97, 138, 63, 110, 310, 101,
        69, 179, 185, 118, 171, 166, 144, 97, 168, 68, 49,
        68, 245, 184, 202, 137, 85, 131, 283, 129, 59, 341,
        87, 65, 102, 265, 276, 252, 90, 100, 55, 61, 92,
        259, 53, 190, 142, 75, 142, 155, 225, 59, 104, 182,
        128, 52, 37, 170, 170, 61, 144, 52, 128, 71, 163,
        150, 97, 160, 178, 48, 270, 202, 111, 85, 42, 170,
        200, 252, 113, 143, 51, 52, 210, 65, 141, 85, 134,
        42, 111, 98, 164, 48, 96, 90, 162, 150, 279, 92,
        83, 128, 102, 302, 198, 95, 53, 134, 144, 232, 81,
        104, 59, 246, 297, 258, 229, 275, 281, 179, 200, 200,
        173, 180, 84, 121, 161, 99, 109, 115, 268, 274, 158,
        107, 83, 103, 272, 85, 280, 336, 281, 118, 317, 235,
        60, 174, 259, 178, 128, 96, 126, 288, 88, 292, 71,
        197, 186, 25, 84, 96, 395, 53, 217, 172, 151, 234,
        59, 70, 220, 268, 152, 47, 74, 295, 101, 151, 127,
        237, 225, 81, 151, 107, 64, 138, 189, 265, 101, 137,
        143, 141, 79, 292, 178, 91, 116, 86, 122, 72, 129,
        142, 90, 158, 39, 196, 222, 277, 89, 196, 202, 155,
        77, 191, 70, 73, 49, 65, 263, 248, 296, 214, 185,
        78, 93, 252, 150, 77, 208, 77, 108, 160, 53, 220,
```

## Drop columns that won't be used for training

```

In [20]: X
Out[20]:
```

	age	bmi	bp	s1	s4	s5	s6
0	0.038076	0.061696	0.021872	-0.044223	-0.002592	0.019908	-0.017646
1	-0.001882	-0.051474	-0.026328	-0.008449	-0.039493	-0.068330	-0.092204
2	0.085299	0.044451	-0.005671	-0.045599	-0.002592	0.002864	-0.025930
3	-0.089063	-0.011595	-0.036656	0.012191	0.034309	0.022692	-0.009362
4	0.005383	-0.036385	0.021872	0.003935	-0.002592	-0.031991	-0.046641
...	...	...	...	...	...	...	...
437	0.041708	0.019662	0.059744	-0.005697	-0.002592	0.031193	0.007207
438	-0.005515	-0.015906	-0.067642	0.049341	0.034309	-0.018118	0.044485
439	0.041708	-0.015906	0.017282	-0.037344	-0.011080	-0.048779	0.015491
440	-0.045472	0.039062	0.001215	0.016318	0.026560	0.044528	-0.025930
441	-0.045472	-0.073030	-0.081414	0.083740	-0.039493	-0.004220	0.003064

442 rows x 7 columns

```
X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

```
((353, 7), (89, 7), (353,), (89,))
```

## Create a linear regression model

```
In [24]: lr.fit(X_train, y_train)

Out[24]: LinearRegression()
```

## Compare the first ten predictions to actual values

```
In [25]: # Make predictions on the test set.

# View examples of the predictions compared to actual disease progression.
```

```
lr_pred = lr.predict(X_test)
```

```
lr_pred[:10]
```

```

181.2071777 , 251.03066105, 105.76063771, 194.56077313,
144.43887399, 229.1423256 )

In [28]: y_test[:10]

Out[28]:
362    321.0
249    215.0
271    127.0
435     64.0
400    175.0
403    275.0
12     179.0
399    232.0
198    142.0
205     99.0
Name: target, dtype: float64

In [29]: X_test

Out[29]:
      age      bmi      bp      s1      s4      s5      s6
362  0.019913  0.104809  0.070073 -0.035968 -0.002592  0.003712  0.040343
249 -0.012780  0.060618  0.052858  0.047965  0.034309  0.070211  0.007207
271  0.038076  0.008883  0.042530 -0.042848 -0.002592 -0.018118  0.007207
435 -0.012780 -0.023451 -0.040099 -0.016704 -0.002592 -0.038459 -0.038357
400 -0.023677  0.045529  0.090730 -0.018080 -0.039493 -0.034524 -0.009362
...     ...     ...     ...     ...     ...     ...
381 -0.070900 -0.089197 -0.074528 -0.042848 -0.002592 -0.012908 -0.054925
213  0.001751 -0.070875 -0.022885 -0.001569 -0.039493 -0.022512  0.007207
134 -0.074533  0.043373 -0.033214  0.012191 -0.039493 -0.027129 -0.046641
49  -0.041840  0.014272 -0.005671 -0.012577  0.071210  0.035462 -0.013504
52  -0.052738 -0.009439 -0.005671  0.039710 -0.002592 -0.018118 -0.013504

89 rows x 7 columns

In [30]: table = X_test.copy()

In [31]: table.shape

Out[31]: (89, 7)

In [32]: table["True Value"] = y_test.copy()

In [33]: table

Out[33]:
      age      bmi      bp      s1      s4      s5      s6  True Value
362  0.019913  0.104809  0.070073 -0.035968 -0.002592  0.003712  0.040343    321.0
249 -0.012780  0.060618  0.052858  0.047965  0.034309  0.070211  0.007207    215.0
271  0.038076  0.008883  0.042530 -0.042848 -0.002592 -0.018118  0.007207    127.0
435 -0.012780 -0.023451 -0.040099 -0.016704 -0.002592 -0.038459 -0.038357     64.0
400 -0.023677  0.045529  0.090730 -0.018080 -0.039493 -0.034524 -0.009362    175.0
...     ...     ...     ...     ...     ...     ...     ...
381 -0.070900 -0.089197 -0.074528 -0.042848 -0.002592 -0.012908 -0.054925    104.0
213  0.001751 -0.070875 -0.022885 -0.001569 -0.039493 -0.022512  0.007207     49.0

```

49	-0.041840	0.014272	-0.005671	-0.012577	0.071210	0.035462	-0.013504	142.0
52	-0.052738	-0.009439	-0.005671	0.039710	-0.002592	-0.018118	-0.013504	59.0

```

In [34]: table["Predicted"] = np.round(lr_pred,2)

In [35]: table

Out[35]:
   age      bmi      bp      s1      s4      s5      s6  True Value  Predicted
362  -0.019913  0.104809  0.070073 -0.035968 -0.002592  0.003712  0.040343      321.0      246.52
249  -0.012780  0.060618  0.052858  0.047965  0.034309  0.070211  0.007207      215.0      237.50
271  0.038076  0.008883  0.042530 -0.042848 -0.002592 -0.018118  0.007207      127.0      168.15
435  -0.012780 -0.023451 -0.040099 -0.016704 -0.002592 -0.038459 -0.038357      64.0      108.45
400  -0.023677  0.045529  0.090730 -0.018080 -0.039493 -0.034524 -0.009362      175.0      181.21
...      ...      ...      ...      ...      ...      ...      ...
381  -0.070900 -0.089197 -0.074528 -0.042848 -0.002592 -0.012908 -0.054925      104.0      86.30
213  0.001751 -0.070875 -0.022885 -0.001569 -0.039493 -0.022512  0.007207      49.0      82.69
134  -0.074533  0.043373 -0.033214  0.012191 -0.039493 -0.027129 -0.046641      103.0      145.36
49   -0.041840  0.014272 -0.005671 -0.012577  0.071210  0.035462 -0.013504      142.0      199.57
52   -0.052738 -0.009439 -0.005671  0.039710 -0.002592 -0.018118 -0.013504      59.0      123.23

89 rows x 9 columns

Calculate the error between predicted and actual values

In [36]: # Print the mean squared error (MSE) for the model's predictions on the test set.
mean_squared_error(y_test,lr_pred)

Out[36]: 3531.2250809792145

Plot lines of best fit for four features

In [37]: # Use Seaborn to create subplots for the four features that have the strongest correlation with the label.
# Also plot a line of best fit for each feature.

line_color = ['color': 'red']
fig, ax = plt.subplots(2,2, figsize=(20,20))

#BMI
ax1 = sns.regplot(x=x_test.bmi, y=lr_pred, line_kws=line_color, ax=ax[0,0])
ax1.set_xlabel("BMI")
ax1.set_ylabel("BMI")


#s5
ax2 = sns.regplot(x=x_test.s5, y=lr_pred, line_kws=line_color, ax=ax[0,1])
ax2.set_xlabel("s5")
ax2.set_ylabel("s5")

#bp
ax3 = sns.regplot(x=x_test.bp, y=lr_pred, line_kws=line_color, ax=ax[1,0])
ax3.set_xlabel("Label")
ax3.set_ylabel("BP")

#s4
ax4 = sns.regplot(x=x_test.s4, y=lr_pred, line_kws=line_color, ax=ax[1,1])
ax4.set_xlabel("Label")
ax4.set_ylabel("s4")

plt.show()

```



In [ ]:

```
fig = plt.figure(figsize=(10, 6))
fig.savefig('MathJaxJax/output/CommonHTML/fonts/TeX/fontdata.js')
```



