

In [1]:

```

1 import numpy as np
2 from sklearn.datasets import load_boston
3 boston = load_boston()
4 import pandas as pd
5 X=pd.DataFrame(boston.data[:100,:], columns=boston.feature_names)
6 y=boston.target[:100]
7 # limit 100 data

```

In [2]:

```
1 X.head()
```

Out[2]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LST
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5

In [3]:

```
1 X.describe()
```

Out[3]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	L
count	100.000000	100.000000	100.000000	100.0	100.000000	100.000000	100.000000	100.000000
mean	0.283301	12.385000	6.668700	0.0	0.471321	6.234410	56.252000	5.177200
std	0.389102	22.966802	3.162352	0.0	0.046836	0.490838	27.476787	1.412000
min	0.006320	0.000000	0.740000	0.0	0.398000	5.399000	2.900000	3.092100
25%	0.050395	0.000000	4.490000	0.0	0.437000	5.926250	35.225000	4.042200
50%	0.102405	0.000000	6.910000	0.0	0.453000	6.130500	57.300000	4.990000
75%	0.221050	18.750000	8.140000	0.0	0.524000	6.433000	81.775000	6.067100
max	1.612820	100.000000	15.040000	0.0	0.538000	8.069000	100.000000	9.222900

In [4]:

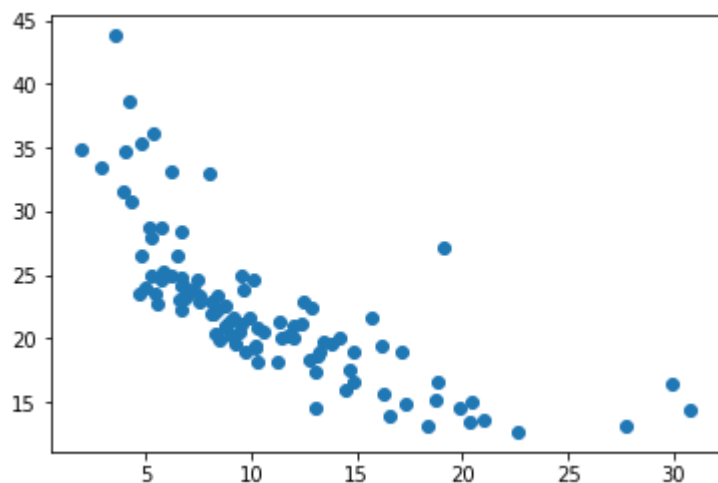
```
1 x = X['LSTAT'].values
```

In [5]:

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 %matplotlib inline
4 plt.scatter(x,y)
```

Out[5]:

<matplotlib.collections.PathCollection at 0x112ca0518>



In [6]:

```
1 from sklearn.linear_model import LinearRegression
2 lin_1d = LinearRegression()
```

In [7]:

```
1 lin_1d.fit(x[:,None], y)
```

Out[7]:

LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=False)

In [8]:

```
1 x.ndim
```

Out[8]:

1

In [9]:

```
1 x[:, np.newaxis].ndim
```

Out[9]:

2

In [10]:

```
1 print(lin_1d.predict(2))
```

[29.09604558]

In [11]:

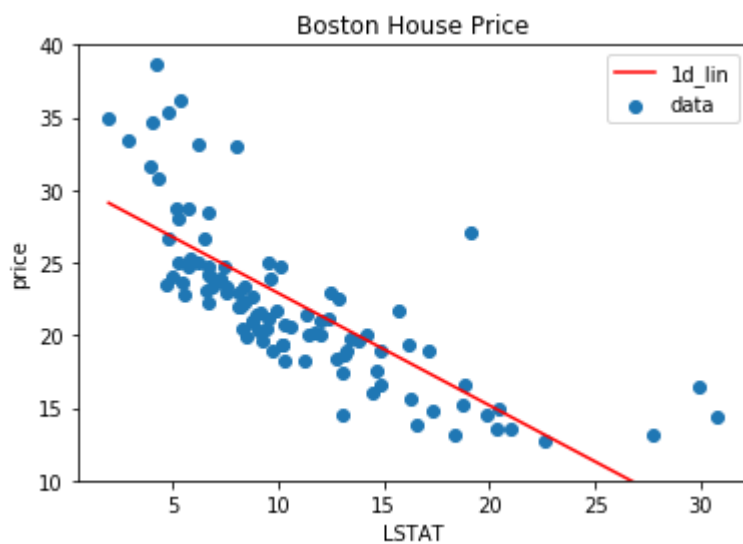
```

1  n = np.linspace(np.min(x), np.max(x), 1000)
2  y_1d_fit = lin_1d.predict(n[:, np.newaxis])
3
4  plt.title("Boston House Price")
5  plt.scatter(x,y,label='data')
6  plt.plot(n, y_1d_fit, 'r', label='1d_lin')
7  plt.ylim(10, 40)
8  plt.xlabel("LSTAT")
9  plt.ylabel("price")
10 plt.legend()
11 plt.plot()

```

Out[11]:

[]



In [12]:

```

1  from sklearn.preprocessing import PolynomialFeatures
2  degree_2 = PolynomialFeatures(degree=2)

```

In [13]:

```
1  degree_2
```

Out[13]:

PolynomialFeatures(degree=2, include_bias=True, interaction_only=False)

In [14]:

```
1  x_2 = degree_2.fit_transform(x[:,None])
```

In [15]:

```

1  x_2
2  lin_2d = LinearRegression()
3  lin_2d.fit(x_2, y)

```

Out[15]:

LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=False)

In [16]:

```

1  n = np.linspace(np.min(x), np.max(x), 1000)
2  y_2d_fit=lin_2d.predict(degree_2.fit_transform(n[:, np.newaxis]))
3  y_1d_fit=lin_1d.predict(n[:, np.newaxis])
4
5  plt.title("Boston House Price")
6  plt.scatter(x,y,label='data')
7  plt.plot(n, y_2d_fit, 'g', label='2d_lin')
8  plt.plot(n, y_1d_fit, 'r', label='1d_lin')
9  plt.ylim(10, 40)
10 plt.xlabel("LSTAT")
11 plt.ylabel("price")
12 plt.legend()
13 plt.plot()

```

Out[16]:

[]



In [17]:

```

1  from sklearn.metrics import mean_squared_error
2
3  mean_squared_error(y, lin_1d.predict(x[:, np.newaxis]))

```

Out[17]:

15.537906007479371

In [18]:

```

1  mean_squared_error(y, lin_2d.predict(x_2))

```

Out[18]:

10.920226905348915

課題 1

まず今回テキストではLSTATのデータ100個に対し、1次、2次関数をフィッティングしました。もしさらに高次の3次関数、4次関数をフィッティングするとどうなるのでしょうか。テキストで学んだscikit learnを使った手法で、3次、4次関数についても、同じように係数を求め、結果をデータ点、1次関数、2次関数、3次関数、4次

関数をまとめてプロットしてみましょう。また学習後の2乗和誤差を3次関数、4次関数についても求め、2乗和誤差が最も小さい関数はどれになるか調べてください。（プロットを見やすくするためデータ点はテキストの通り100個で行ってください。）

回答：

ソースコードとプロット図とそれぞれの関数の二乗誤差は下記にあり、最も小さいのは**4次関数**になる。

In [19]:

```

1 degree_3 = PolynomialFeatures(degree=3)
2 degree_4 = PolynomialFeatures(degree=4)
3
4 x_3 = degree_3.fit_transform(x[:,None])
5 x_4 = degree_4.fit_transform(x[:,None])
6
7 lin_3d = LinearRegression()
8 lin_4d = LinearRegression()
9
10 lin_3d.fit(x_3, y)
11 lin_4d.fit(x_4, y)
12
13 n = np.linspace(np.min(x), np.max(x), 1000)
14
15 y_3d_fit=lin_3d.predict(degree_3.fit_transform(n[:, np.newaxis]))
16 y_4d_fit=lin_4d.predict(degree_4.fit_transform(n[:, np.newaxis]))
17
18 plt.title("Boston House Price")
19 plt.scatter(x,y,label='data')
20 plt.plot(n, y_2d_fit, 'g', label='2d_lin')
21 plt.plot(n, y_1d_fit, 'r', label='1d_lin')
22 plt.plot(n, y_3d_fit, 'b', label='3d_lin')
23 plt.plot(n, y_4d_fit, 'y', label='4d_lin')
24 plt.ylim(10, 40)
25 plt.xlabel("LSTAT")
26 plt.ylabel("price")
27 plt.legend()
28 plt.plot()
29
30 print(mean_squared_error(y, lin_1d.predict(x[:,np.newaxis])))
31 print(mean_squared_error(y, lin_2d.predict(x_2)))
32 print(mean_squared_error(y, lin_3d.predict(x_3)))
33 print(mean_squared_error(y, lin_4d.predict(x_4)))

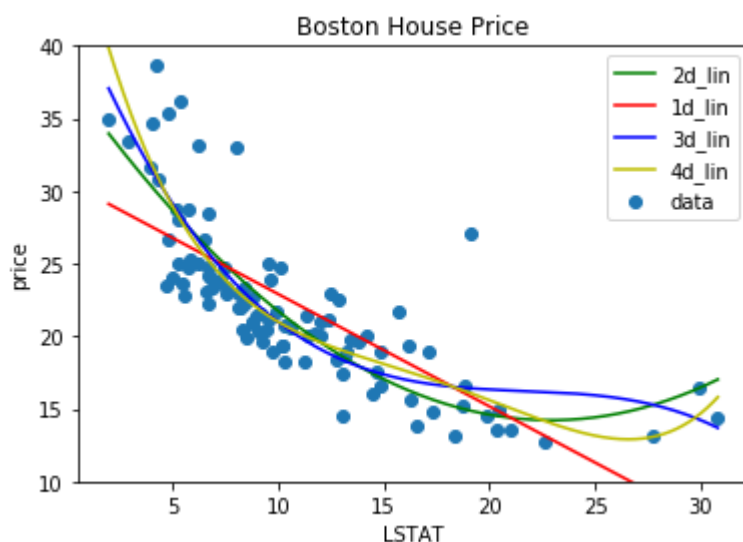
```

15.537906007479371

10.920226905348915

10.05522570974638

9.617390715416688



課題2

もう少し回帰の練習をつむため、今度は同じボストンの住宅価格のデータセットにおいて特徴量としてLSTATではなくAGEを使ってみましょう。

```
x=X['AGE'].values
```

とすればデータを読み込むことができます。

このデータについてもテキストで学んだscikit learnを使った手法で1次から4次関数についてまで、同じように係数を求め、結果をデータ点、1次関数、2次関数、3次関数、4次関数をまとめてプロットしてみましょう。また学習後の2乗和誤差をそれぞれについても求め、2乗和誤差が最も小さい関数はどれになるか調べてください。

In [23]:

```

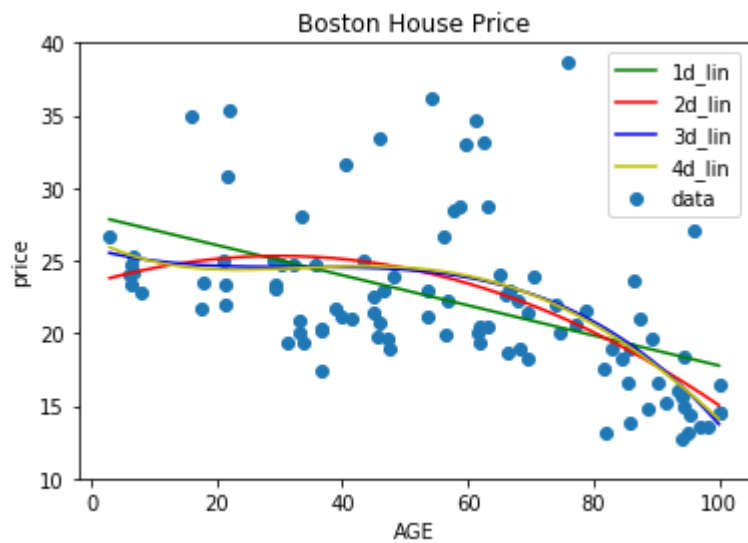
1  # モジュールインポート以外をもう一度最初からやってみる
2  boston = load_boston()
3  X = pd.DataFrame(boston.data[:100,:], columns=boston.feature_names)
4  y=boston.target[:100]
5  #データを100個に限っているのは後でグラフをプロットしたときの様子を見やすくするため
6  x = X['AGE'].values
7
8  y_pred = []
9
10 for i in [1, 2, 3, 4]:
11     # 線形回帰モデルを使う
12     lin_ = LinearRegression()
13     if i == 1:
14         x_ = x[:, np.newaxis]
15     else:
16         degree_ = PolynomialFeatures(degree=i)
17         x_ = degree_.fit_transform(x[:, None])
18
19     lin_.fit(x_, y)
20     n = np.linspace(np.min(x), np.max(x), 1000)
21
22     if i==1:
23         y_pred.append(lin_.predict(n[:, np.newaxis]))
24         print('1次関数の2乗誤差は')
25         print(mean_squared_error(y, lin_.predict(x_)))
26     else:
27         y_pred.append(lin_.predict(degree_.fit_transform(n[:, np.newaxis])))
28         print('{}次関数の2乗誤差は'.format(i))
29         print(mean_squared_error(y, lin_.predict(x_)))
30
31
32 plt.title("Boston House Price")
33 plt.scatter(x,y,label='data')
34 plt.plot(n, y_pred[0], 'g', label='1d_lin')
35 plt.plot(n, y_pred[1], 'r', label='2d_lin')
36 plt.plot(n, y_pred[2], 'b', label='3d_lin')
37 plt.plot(n, y_pred[3], 'y', label='4d_lin')
38 plt.ylim(10, 40)
39 plt.xlabel("AGE")
40 plt.ylabel("price")
41 plt.legend()
42 plt.plot()

```

1次関数の2乗誤差は
26.754508704139507
2次関数の2乗誤差は
24.39708813651596
3次関数の2乗誤差は
23.97824219967744
4次関数の2乗誤差は
23.956816976819297

Out[23]:

[]



回答

上記ソースコードでプロットしており、2乗誤差が最小なのは**4次関数**となる。