

assignment2

May 10, 2025

B. Practical Overview of the steps 1. Load the data and get an overview of the data 2. Perform simple linear regressions 3. Use the simple linear regression models 4. Perform multiple linear regressions 5. Use the multiple linear regression model Steps in detail Load the data and get an overview of the data Load the data file Boston.rda or Boston.csv .

0.0.1 Steps in Detail

Load the Data and Get an Overview of the Data

- Load the data file Boston.csv.

```
[2]: import pandas as pd
df = pd.read_csv("dataset/Boston.csv", index_col=0)
```

0.0.2 Display the Number of Predictors and Their Names

```
[3]: print(df.shape[1])
print(df.columns.tolist())
```

14

```
['crim', 'zn', 'indus', 'chas', 'nox', 'rm', 'age', 'dis', 'rad', 'tax',
'ptratio', 'black', 'lstat', 'medv']
```

0.0.3 Print a Statistical Summary of the Predictors and the Response medv

```
[4]: print(df.describe())
```

| | crim | zn | indus | chas | nox | rm | \ |
|-------|------------|------------|------------|------------|------------|------------|---|
| count | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 | |
| mean | 3.613524 | 11.363636 | 11.136779 | 0.069170 | 0.554695 | 6.284634 | |
| std | 8.601545 | 23.322453 | 6.860353 | 0.253994 | 0.115878 | 0.702617 | |
| min | 0.006320 | 0.000000 | 0.460000 | 0.000000 | 0.385000 | 3.561000 | |
| 25% | 0.082045 | 0.000000 | 5.190000 | 0.000000 | 0.449000 | 5.885500 | |
| 50% | 0.256510 | 0.000000 | 9.690000 | 0.000000 | 0.538000 | 6.208500 | |
| 75% | 3.677083 | 12.500000 | 18.100000 | 0.000000 | 0.624000 | 6.623500 | |
| max | 88.976200 | 100.000000 | 27.740000 | 1.000000 | 0.871000 | 8.780000 | |

| | age | dis | rad | tax | ptratio | black | \ |
|-------|------------|------------|------------|------------|------------|------------|---|
| count | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 | 506.000000 | |
| mean | 68.574901 | 3.795043 | 9.549407 | 408.237154 | 18.455534 | 356.674032 | |

| | | | | | | |
|-----|------------|-----------|-----------|------------|-----------|------------|
| std | 28.148861 | 2.105710 | 8.707259 | 168.537116 | 2.164946 | 91.294864 |
| min | 2.900000 | 1.129600 | 1.000000 | 187.000000 | 12.600000 | 0.320000 |
| 25% | 45.025000 | 2.100175 | 4.000000 | 279.000000 | 17.400000 | 375.377500 |
| 50% | 77.500000 | 3.207450 | 5.000000 | 330.000000 | 19.050000 | 391.440000 |
| 75% | 94.075000 | 5.188425 | 24.000000 | 666.000000 | 20.200000 | 396.225000 |
| max | 100.000000 | 12.126500 | 24.000000 | 711.000000 | 22.000000 | 396.900000 |

| | | |
|-------|------------|------------|
| | lstat | medv |
| count | 506.000000 | 506.000000 |
| mean | 12.653063 | 22.532806 |
| std | 7.141062 | 9.197104 |
| min | 1.730000 | 5.000000 |
| 25% | 6.950000 | 17.025000 |
| 50% | 11.360000 | 21.200000 |
| 75% | 16.955000 | 25.000000 |
| max | 37.970000 | 50.000000 |

0.0.4 Display the number of data points:

```
[5]: print(df.shape[0])
```

506

0.0.5 Display the data in a table (subset of rows is sufficient):

```
[6]: print(df.head())
```

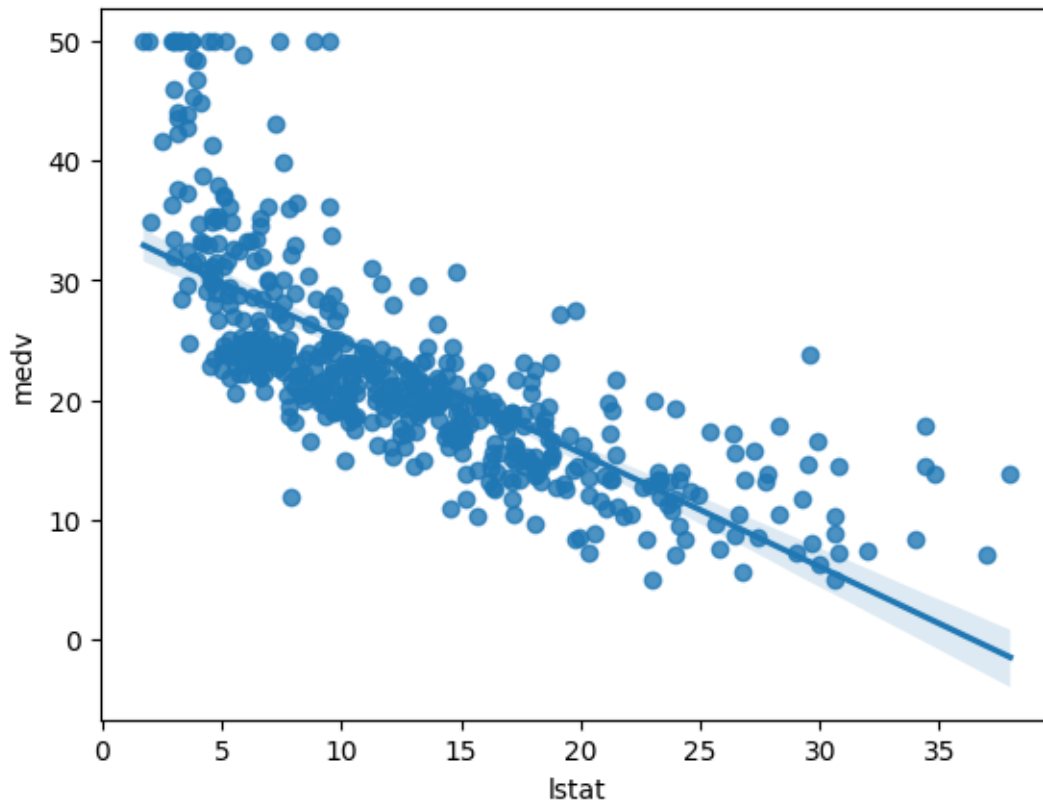
| | | | | | | | | | | | | |
|---|---------|------|-------|------|-------|-------|------|--------|-----|-----|---------|---|
| | crim | zn | indus | chas | nox | rm | age | dis | rad | tax | ptratio | \ |
| 1 | 0.00632 | 18.0 | 2.31 | 0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1 | 296 | 15.3 | |
| 2 | 0.02731 | 0.0 | 7.07 | 0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2 | 242 | 17.8 | |
| 3 | 0.02729 | 0.0 | 7.07 | 0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2 | 242 | 17.8 | |
| 4 | 0.03237 | 0.0 | 2.18 | 0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3 | 222 | 18.7 | |
| 5 | 0.06905 | 0.0 | 2.18 | 0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3 | 222 | 18.7 | |

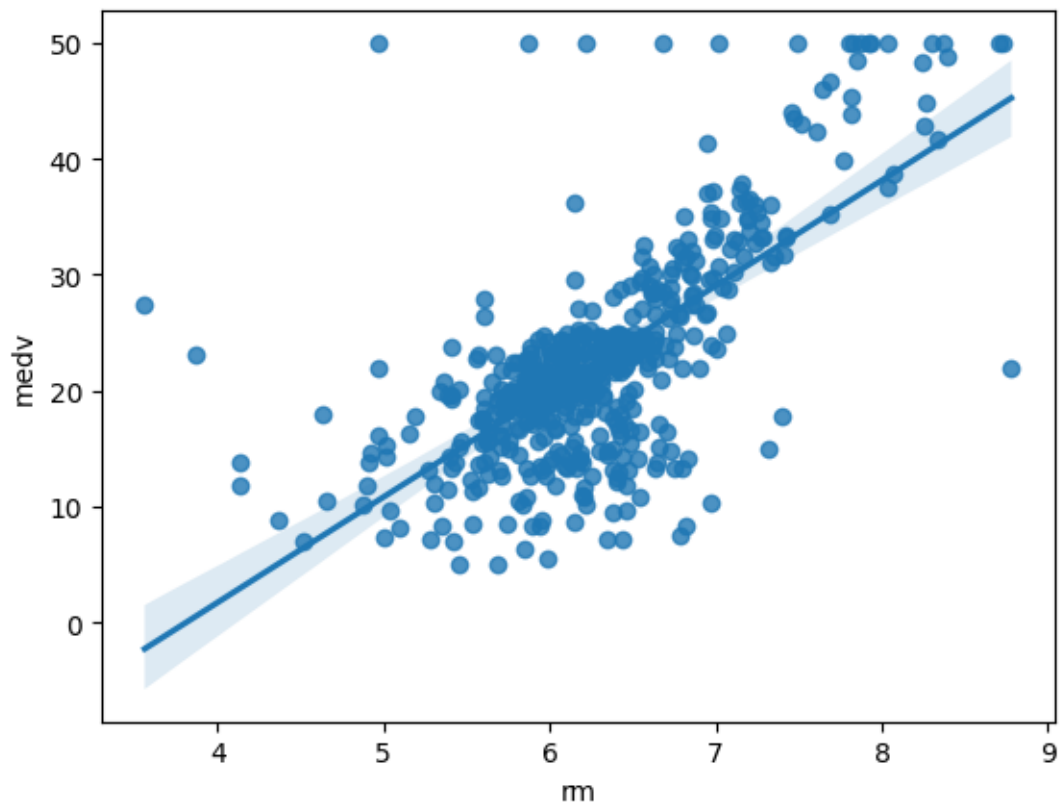
| | | | |
|---|--------|-------|------|
| | black | lstat | medv |
| 1 | 396.90 | 4.98 | 24.0 |
| 2 | 396.90 | 9.14 | 21.6 |
| 3 | 392.83 | 4.03 | 34.7 |
| 4 | 394.63 | 2.94 | 33.4 |
| 5 | 396.90 | 5.33 | 36.2 |

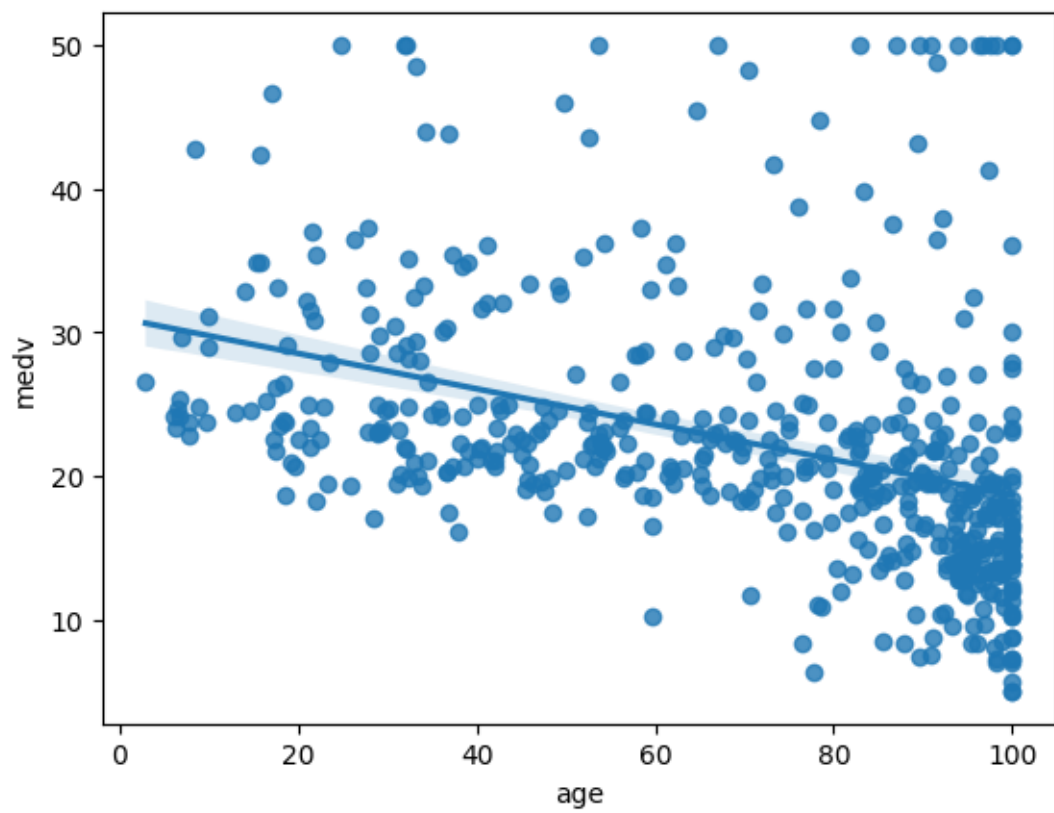
0.0.6 Plot some predictors (at least two) against the response values. We choose lstat, rm, and age. I added crim and pupol-teacher ratio

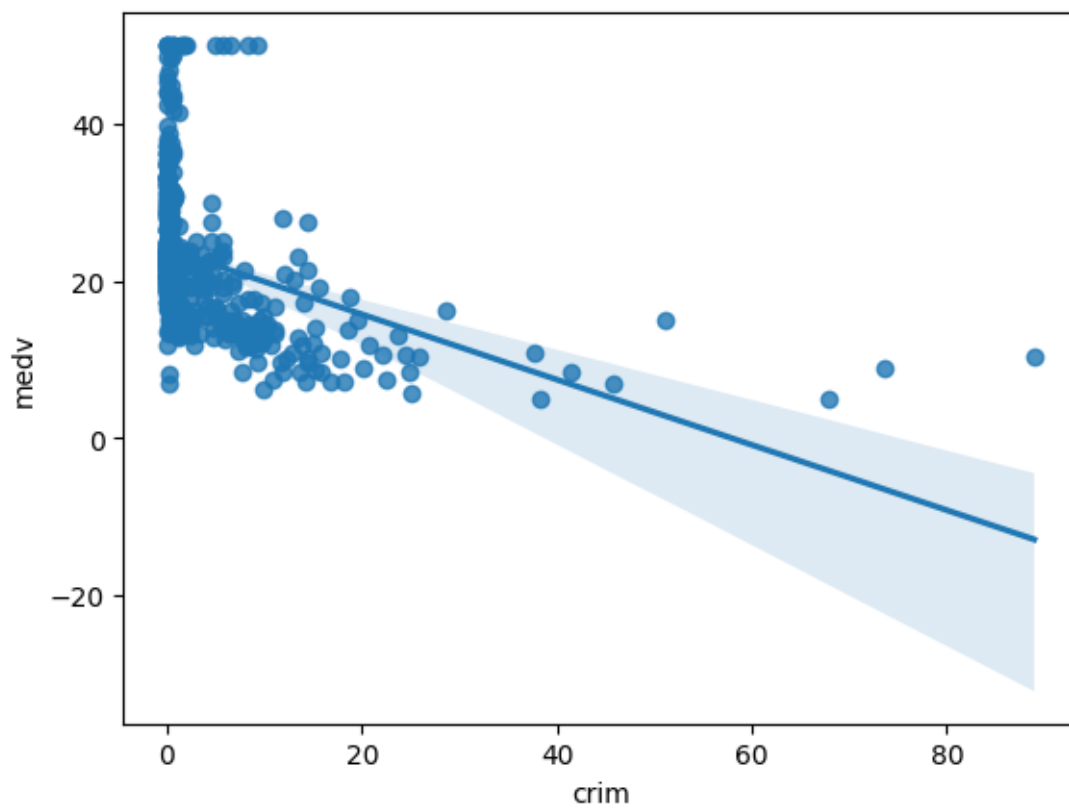
```
[7]: import seaborn as sns
import matplotlib.pyplot as plt
sns.regplot(x='lstat', y='medv', data=df)
plt.show()
```

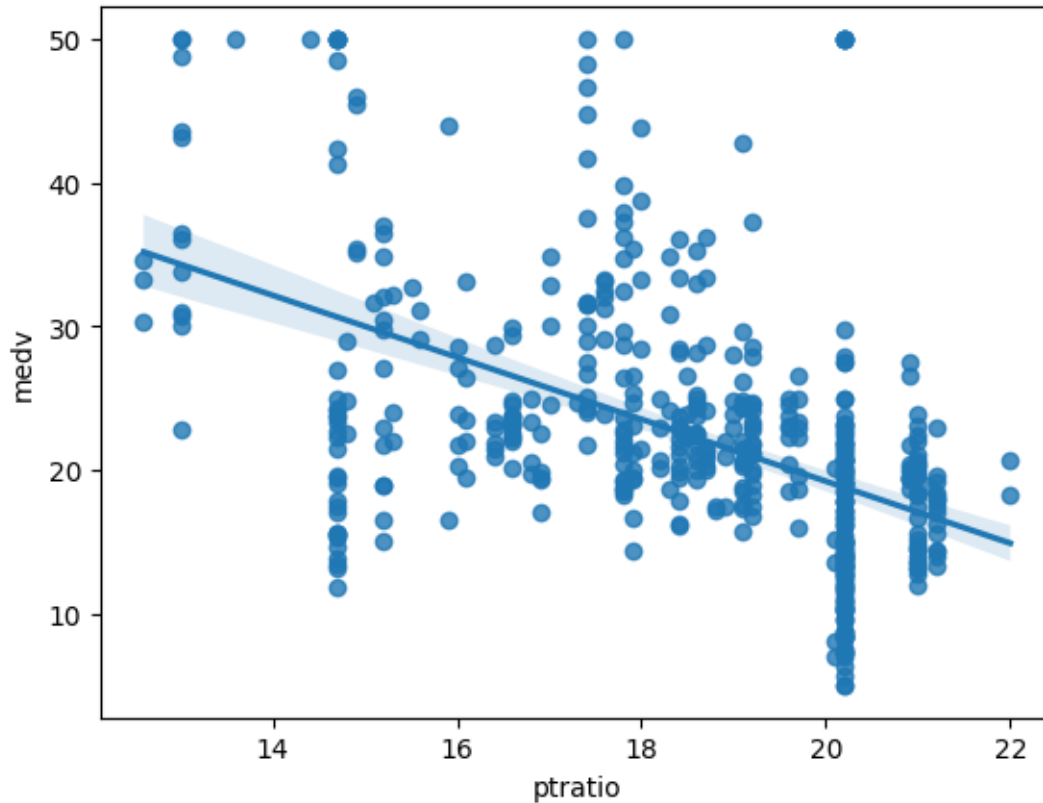
```
sns.regplot(x='rm', y='medv', data=df)
plt.show()
sns.regplot(x='age', y='medv', data=df)
plt.show()
sns.regplot(x='crim', y='medv', data=df)
plt.show()
sns.regplot(x='ptratio', y='medv', data=df)
plt.show()
```











0.0.7 Perform Simple Linear Regressions

Fit a simple linear regression model with `medv` as the response and some (at least two) predictors individually. We choose the following predictors:

- `lstat`
- `rm`
- `age`

i added

- `crim`
- `ptratio`

```
[8]: import statsmodels.api as sm
print("this stats for lstat")
X = sm.add_constant(df['lstat'])
model_lstat = sm.OLS(df['medv'], X).fit()
print(model_lstat.summary())
```

this stats for lstat

OLS Regression Results

=====

```

Dep. Variable:          medv    R-squared:          0.544
Model:                  OLS     Adj. R-squared:       0.543
Method:                 Least Squares    F-statistic:        601.6
Date:                   Sat, 10 May 2025    Prob (F-statistic):  5.08e-88
Time:                   16:19:53    Log-Likelihood:     -1641.5
No. Observations:      506    AIC:                3287.
Df Residuals:          504    BIC:                3295.
Df Model:               1
Covariance Type:       nonrobust

```

```

=====
              coef    std err          t      P>|t|      [0.025      0.975]
-----
const         34.5538     0.563     61.415     0.000     33.448     35.659
lstat        -0.9500     0.039    -24.528     0.000     -1.026     -0.874
=====

Omnibus:                 137.043    Durbin-Watson:           0.892
Prob(Omnibus):            0.000    Jarque-Bera (JB):        291.373
Skew:                     1.453    Prob(JB):                 5.36e-64
Kurtosis:                  5.319    Cond. No.                  29.7
=====

```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```

[9]: print("this stats for rm")
      X = sm.add_constant(df['rm'])
      model_rm = sm.OLS(df['medv'], X).fit()
      print(model_rm.summary())

```

this stats for rm

```

                                OLS Regression Results
=====
Dep. Variable:          medv    R-squared:          0.484
Model:                  OLS     Adj. R-squared:       0.483
Method:                 Least Squares    F-statistic:        471.8
Date:                   Sat, 10 May 2025    Prob (F-statistic):  2.49e-74
Time:                   16:19:53    Log-Likelihood:     -1673.1
No. Observations:      506    AIC:                3350.
Df Residuals:          504    BIC:                3359.
Df Model:               1
Covariance Type:       nonrobust

```

```

=====
              coef    std err          t      P>|t|      [0.025      0.975]
-----
const        -34.6706     2.650    -13.084     0.000    -39.877    -29.465
rm             9.1021     0.419     21.722     0.000      8.279      9.925
=====

```


| | | | |
|----------------|---------|-------------------|-----------|
| Omnibus: | 102.585 | Durbin-Watson: | 0.684 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 612.449 |
| Skew: | 0.726 | Prob(JB): | 1.02e-133 |
| Kurtosis: | 8.190 | Cond. No. | 58.4 |

=====

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
[10]: print("this stats for age")
      X = sm.add_constant(df['age'])
      model_age = sm.OLS(df['medv'], X).fit()
      print(model_age.summary())
```

this stats for age

OLS Regression Results

=====

| | | | |
|-------------------|------------------|---------------------|----------|
| Dep. Variable: | medv | R-squared: | 0.142 |
| Model: | OLS | Adj. R-squared: | 0.140 |
| Method: | Least Squares | F-statistic: | 83.48 |
| Date: | Sat, 10 May 2025 | Prob (F-statistic): | 1.57e-18 |
| Time: | 16:19:53 | Log-Likelihood: | -1801.5 |
| No. Observations: | 506 | AIC: | 3607. |
| Df Residuals: | 504 | BIC: | 3615. |
| Df Model: | 1 | | |
| Covariance Type: | nonrobust | | |

=====

| | coef | std err | t | P> t | [0.025 | 0.975] |
|-------|---------|---------|--------|-------|--------|--------|
| const | 30.9787 | 0.999 | 31.006 | 0.000 | 29.016 | 32.942 |
| age | -0.1232 | 0.013 | -9.137 | 0.000 | -0.150 | -0.097 |

=====

| | | | |
|----------------|---------|-------------------|-----------|
| Omnibus: | 170.034 | Durbin-Watson: | 0.613 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 456.983 |
| Skew: | 1.671 | Prob(JB): | 5.85e-100 |
| Kurtosis: | 6.240 | Cond. No. | 195. |

=====

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
[11]: print("this stats for crim")
      X = sm.add_constant(df['crim'])
      model_crim = sm.OLS(df['medv'], X).fit()
      print(model_crim.summary())
```

this stats for crim

OLS Regression Results

```

=====
Dep. Variable:          medv    R-squared:                0.151
Model:                  OLS    Adj. R-squared:             0.149
Method:                 Least Squares    F-statistic:        89.49
Date:                   Sat, 10 May 2025    Prob (F-statistic):    1.17e-19
Time:                   16:19:53    Log-Likelihood:        -1798.9
No. Observations:      506    AIC:                  3602.
Df Residuals:          504    BIC:                  3610.
Df Model:               1
Covariance Type:       nonrobust
=====

```

```

=====
              coef    std err          t      P>|t|      [0.025    0.975]
-----
const         24.0331      0.409     58.740      0.000      23.229      24.837
crim          -0.4152      0.044     -9.460      0.000      -0.501      -0.329
=====

```

```

=====
Omnibus:            139.832    Durbin-Watson:           0.713
Prob(Omnibus):      0.000    Jarque-Bera (JB):        295.404
Skew:               1.490    Prob(JB):                7.14e-65
Kurtosis:           5.264    Cond. No.                10.1
=====

```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```

[12]: print("this stats for ptratio")
X = sm.add_constant(df['ptratio'])
model_ptratio = sm.OLS(df['medv'], X).fit()
print(model_ptratio.summary())

```

this stats for ptratio

OLS Regression Results

```

=====
Dep. Variable:          medv    R-squared:                0.258
Model:                  OLS    Adj. R-squared:             0.256
Method:                 Least Squares    F-statistic:        175.1
Date:                   Sat, 10 May 2025    Prob (F-statistic):    1.61e-34
Time:                   16:19:53    Log-Likelihood:        -1764.8
No. Observations:      506    AIC:                  3534.
Df Residuals:          504    BIC:                  3542.
Df Model:               1
Covariance Type:       nonrobust
=====

```

```

=====
              coef    std err          t      P>|t|      [0.025    0.975]
-----
const         62.3446      3.029     20.581      0.000      56.393      68.296

```

| | | | | | | |
|----------------|---------|--------|-------------------|-------|--------|----------|
| pvalue | -2.1572 | 0.163 | -13.233 | 0.000 | -2.477 | -1.837 |
| ===== | | | | | | |
| Omnibus: | | 92.924 | Durbin-Watson: | | | 0.725 |
| Prob(Omnibus): | | 0.000 | Jarque-Bera (JB): | | | 191.444 |
| Skew: | | 1.001 | Prob(JB): | | | 2.68e-42 |
| Kurtosis: | | 5.252 | Cond. No. | | | 160. |
| ===== | | | | | | |

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

0.0.8 Percent of low socioeconomic status households (lstat) vs house price

From the summary stats, `lstat` ranges from ~1.7 to ~38, showing a wide spread and likely right-skewed distribution.

When we plot it, we see a clear **negative relationship** — as `lstat` increases, `medv` decreases. This makes sense: poorer neighborhoods tend to have lower house prices, and the data strongly support that.

From the regression, each +1% in low-status households is linked to about a **\$950 drop** in median price. It's the strongest predictor so far, explaining **54% of the variation** in prices.

In simple terms, wealthier neighborhoods have much higher home values — which is exactly what we'd expect.

0.0.9 Average number of rooms per house (rm) vs house price

Summary stats show `rm` ranges from ~3.6 to ~8.8 rooms, with a fairly balanced spread.

On the plot, we see a **strong positive relationship** — more rooms per house clearly lead to higher prices. Again, this feels intuitive: larger homes tend to be more expensive.

The regression tells us each additional room adds about **\$9,100** to the median price. It explains **48% of the variation**, making it one of the top predictors.

Simply put, bigger houses are worth more — no surprise, but it's great to see the numbers back it up.

0.0.10 Age of houses (age) vs house price

The age variable ranges from ~2.9% to 100% old units, with a right-skewed pattern (lots of old houses).

The plot shows a **negative relationship**, though weaker than `lstat` or `rm`. As houses get older, prices tend to go down slightly. This makes sense: newer homes or neighborhoods often attract higher prices because of better condition or updated amenities.

The regression shows each +1% increase in old units lowers price by about **\$120**. But it only explains **14% of the variation** — so it's a weaker predictor.

In short, older neighborhoods slightly reduce house values, but they're not the main price driver.

0.0.11 Crime rate per capita (**crim**) vs house price

Summary stats show **crim** ranges from ~0.006 to ~89 — a huge range with major right skew and clear outliers.

The plot shows a **negative but noisy relationship** — higher crime is linked to lower prices, but the scatter is messy. This aligns with intuition: people prefer safer neighborhoods, but the relationship here isn't as tight as we might expect.

Regression tells us each unit increase in crime rate reduces median price by about **\$420**, but it explains only **15% of the variation**.

In simple terms, higher crime weakly lowers house prices, but it's not a main factor compared to wealth or house size.

0.0.12 Pupil-teacher ratio (**ptratio**) vs house price

ptratio ranges from ~12.6 to ~22, with moderate spread.

The plot reveals a **negative relationship** — as class sizes get bigger (more students per teacher), home prices drop. This makes sense because good schools are a key selling point for families, and overcrowded schools can signal lower community investment.

The regression shows each additional student per teacher cuts median price by about **\$2,160**, explaining **26% of the variation**.

In short, school quality matters — neighborhoods with better student-teacher ratios tend to have more expensive homes.

0.0.13 Notes

- All predictors have $p < 0.001 \rightarrow$ **statistically significant**.
- **lstat** and **rm** are the top predictors (highest R^2).
- **crim**, **age**, **ptratio** have much smaller explanatory power (low R^2).

0.0.14 Obtain a confidence interval for the coefficient estimates for the individual models.

```
[13]: print("this is confidence interval for lstat")
      print(model_lstat.conf_int())
      print("this is confidence interval for rm")
      print(model_rm.conf_int())
      print("this is confidence interval for age")
```

```
print(model_age.conf_int())
print("this is confidence interval for crim")
print(model_crim.conf_int())
print("this is confidence interval for ptratio")
print(model_ptratio.conf_int())
```

```
this is confidence interval for lstat
      0      1
const 33.448457 35.659225
lstat -1.026148 -0.873951
this is confidence interval for rm
      0      1
const -39.876641 -29.464601
rm      8.278855  9.925363
this is confidence interval for age
      0      1
const 29.015752 32.941604
age    -0.149647 -0.096679
this is confidence interval for crim
      0      1
const 23.229272 24.83694
crim   -0.501421 -0.32896
this is confidence interval for ptratio
      0      1
const 56.393267 68.295988
ptratio -2.477454 -1.836897
```

0.0.15 Interpretation of confidence intervals

Percent of low socioeconomic status households (lstat)

From the confidence interval, the intercept is between 33.45 and 35.66, and the slope for lstat is between **−1.03** and **−0.87**.

What this tells us:

We're 95% confident that for each +1% increase in low-status households, median price goes down by between **\$870** and **\$1,030**.

The interval is fully negative → this is a strong, consistently negative effect.

Average number of rooms (rm)

The intercept ranges between −39.88 and −29.46, and the slope for rm is between **8.28** and **9.93**.

What this tells us:

We're 95% confident that each extra room adds **\$8,280** to **\$9,930** to the median price.

The interval is fully positive → this is a strong and reliably positive effect.

Age of houses (**age**)

The intercept ranges between 29.02 and 32.94,
and the slope for **age** is between **−0.15** and **−0.10**.

What this tells us:

We're 95% confident that every +1% in old units reduces median price by between **\$100** and **\$150**.
The effect is consistently negative, though smaller compared to the top predictors.

Crime rate per capita (**crim**)

The intercept is between 23.23 and 24.84,
and the slope for **crim** is between **−0.50** and **−0.33**.

What this tells us:

We're 95% confident that each unit increase in crime rate lowers the median price by between **\$330** and **\$500**.

The relationship is moderately negative and quite consistent.

Pupil-teacher ratio (**ptratio**)

The intercept is between 56.39 and 68.30,
and the slope for **ptratio** is between **−2.48** and **−1.84**.

What this tells us:

We're 95% confident that each extra student per teacher reduces median price by between **\$1,840** and **\$2,480**.

The effect is clearly negative and meaningful.

0.0.16 Summary notes

- All predictors have **tight, narrow confidence intervals** → we're pretty confident in the estimates.
- None of the intervals cross zero → all effects are statistically significant.
- The direction of the effects matches our earlier findings:
 - **lstat**, **age**, **crim**, **ptratio** → negative effect on price
 - **rm** → positive effect on price

0.0.17 Use the Simple Linear Regression Models

Predict the **medv** response values for some selected predictor values. Calculate the prediction intervals for these values.

```
[14]: new_lstat = sm.add_constant(pd.DataFrame({'lstat': [5, 10, 15]}))  
      print(model_lstat.get_prediction(new_lstat).summary_frame(alpha=0.05))
```

| | mean | mean_se | mean_ci_lower | mean_ci_upper | obs_ci_lower | \ |
|---|-----------|----------|---------------|---------------|--------------|---|
| 0 | 29.803594 | 0.405247 | 29.007412 | 30.599776 | 17.565675 | |
| 1 | 25.053347 | 0.294814 | 24.474132 | 25.632563 | 12.827626 | |
| 2 | 20.303101 | 0.290893 | 19.731588 | 20.874613 | 8.077742 | |

| | obs_ci_upper |
|---|--------------|
| 0 | 42.041513 |
| 1 | 37.279068 |
| 2 | 32.528459 |

```
[15]: new_rm = sm.add_constant(pd.DataFrame({'rm': [5, 6.5, 8]}))
      print(model_rm.get_prediction(new_rm).summary_frame(alpha=0.05))
```

| | mean | mean_se | mean_ci_lower | mean_ci_upper | obs_ci_lower | \ |
|---|-----------|----------|---------------|---------------|--------------|---|
| 0 | 10.839924 | 0.613410 | 9.634769 | 12.045079 | -2.214474 | |
| 1 | 24.493088 | 0.307657 | 23.888639 | 25.097536 | 11.480391 | |
| 2 | 38.146251 | 0.776633 | 36.620414 | 39.672088 | 25.058353 | |

| | obs_ci_upper |
|---|--------------|
| 0 | 23.894322 |
| 1 | 37.505784 |
| 2 | 51.234149 |

```
[16]: new_age = sm.add_constant(pd.DataFrame({'age': [25, 50, 75]}))
      print(model_age.get_prediction(new_age).summary_frame(alpha=0.05))
```

| | mean | mean_se | mean_ci_lower | mean_ci_upper | obs_ci_lower | \ |
|---|-----------|----------|---------------|---------------|--------------|---|
| 0 | 27.899610 | 0.699094 | 26.526112 | 29.273107 | 11.090368 | |
| 1 | 24.820542 | 0.454307 | 23.927973 | 25.713110 | 8.043748 | |
| 2 | 21.741474 | 0.388844 | 20.977518 | 22.505429 | 4.971031 | |

| | obs_ci_upper |
|---|--------------|
| 0 | 44.708852 |
| 1 | 41.597335 |
| 2 | 38.511917 |

```
[17]: new_crim = sm.add_constant(pd.DataFrame({'crim': [0.1, 1, 5]}))
      print(model_crim.get_prediction(new_crim).summary_frame(alpha=0.05))
```

| | mean | mean_se | mean_ci_lower | mean_ci_upper | obs_ci_lower | \ |
|---|-----------|----------|---------------|---------------|--------------|---|
| 0 | 23.991587 | 0.407461 | 23.191056 | 24.792118 | 7.304363 | |
| 1 | 23.617916 | 0.394210 | 22.843418 | 24.392413 | 6.931921 | |
| 2 | 21.957155 | 0.382030 | 21.206588 | 22.707721 | 5.272253 | |

| | obs_ci_upper |
|---|--------------|
| 0 | 40.678811 |
| 1 | 40.303911 |
| 2 | 38.642056 |

```
[18]: new_ptratio = sm.add_constant(pd.DataFrame({'ptratio': [15, 18, 21]}))
      print(model_ptratio.get_prediction(new_ptratio).summary_frame(alpha=0.05))
```

| | mean | mean_se | mean_ci_lower | mean_ci_upper | obs_ci_lower | \ |
|---|-----------|----------|---------------|---------------|--------------|---|
| 0 | 29.986998 | 0.664554 | 28.681360 | 31.292636 | 14.350492 | |
| 1 | 23.515472 | 0.360312 | 22.807574 | 24.223370 | 7.917500 | |
| 2 | 17.043946 | 0.544393 | 15.974387 | 18.113505 | 1.425381 | |

| | obs_ci_upper |
|---|--------------|
| 0 | 45.623504 |
| 1 | 39.113444 |
| 2 | 32.662511 |

For the two variables I added, `crim` and `ptratio`, I chose the intervals after looking at the summary statistics at the beginning of the experiment; what I ended up with is `crim = [0.1, 1, 5]` and `ptratio = [15, 18, 21]` because these cover low, medium, and high values without going into extreme outliers, giving a realistic and interpretable range.

0.0.18 Interpretation of predicted house prices and intervals

lstat (low socioeconomic households)

- 5% → \$29.8k (range \$17.6k–\$42.0k)
- 10% → \$25.1k (range \$12.8k–\$37.3k)
- 15% → \$20.3k (range \$8.1k–\$32.5k)

As `lstat` increases, prices drop steadily. Poorer neighborhoods have lower home values. The prediction intervals shrink slightly as `lstat` rises, suggesting a bit less uncertainty in low-price areas.

rm (number of rooms)

- 5 rooms → \$10.8k (–\$2.2k–\$23.9k)
- 6.5 rooms → \$24.5k (\$11.5k–\$37.5k)
- 8 rooms → \$38.1k (\$25.1k–\$51.2k)

More rooms strongly increase price. The intervals widen as room count increases, showing more price variability among large homes.

age (percent old houses)

- 25% → \$27.9k (\$11.1k–\$44.7k)
- 50% → \$24.8k (\$8.0k–\$41.6k)
- 75% → \$21.7k (\$5.0k–\$38.5k)

Older areas show a slight price drop, but intervals stay fairly wide across all levels. This suggests mixed price outcomes no matter the age.

crim (crime rate)

- 0.1 → \$24.0k (\$7.3k–\$40.7k)
- 1 → \$23.6k (\$6.9k–\$40.3k)
- 5 → \$22.0k (\$5.3k–\$38.6k)

Higher crime slightly lowers price, but the intervals stay wide, meaning the effect of crime on price can vary a lot between neighborhoods.

ptratio (pupil-teacher ratio)

- 15 → \$30.0k (\$14.4k–\$45.6k)
- 18 → \$23.5k (\$7.9k–\$39.1k)
- 21 → \$17.0k (\$1.4k–\$32.7k)

Better schools (lower ratios) raise prices. The intervals widen as the ratio increases, reflecting more price uncertainty in areas with crowded schools.

0.0.19 Perform Multiple Linear Regressions

Fit medv as the response with the previously selected predictors (lstat, rm, age, crim, ptratio) altogether.

```
[19]: X = sm.add_constant(df[['lstat', 'rm', 'age', 'crim', 'ptratio']])
      model_5 = sm.OLS(df['medv'], X).fit()
      print(model_5.summary())
```

| OLS Regression Results | | | | | | |
|------------------------|------------------|---------|---------------------|-----------|--------|--------|
| ===== | | | | | | |
| Dep. Variable: | medv | | R-squared: | 0.683 | | |
| Model: | OLS | | Adj. R-squared: | 0.680 | | |
| Method: | Least Squares | | F-statistic: | 215.9 | | |
| Date: | Sat, 10 May 2025 | | Prob (F-statistic): | 2.24e-122 | | |
| Time: | 16:19:53 | | Log-Likelihood: | -1549.2 | | |
| No. Observations: | 506 | | AIC: | 3110. | | |
| Df Residuals: | 500 | | BIC: | 3136. | | |
| Df Model: | 5 | | | | | |
| Covariance Type: | nonrobust | | | | | |
| ===== | | | | | | |
| | coef | std err | t | P> t | [0.025 | 0.975] |
| ----- | | | | | | |
| const | 17.5695 | 3.985 | 4.409 | 0.000 | 9.741 | 25.398 |
| lstat | -0.5836 | 0.054 | -10.889 | 0.000 | -0.689 | -0.478 |
| rm | 4.4604 | 0.436 | 10.235 | 0.000 | 3.604 | 5.317 |
| age | 0.0185 | 0.011 | 1.746 | 0.082 | -0.002 | 0.039 |
| crim | -0.0699 | 0.031 | -2.264 | 0.024 | -0.130 | -0.009 |
| ptratio | -0.9049 | 0.119 | -7.610 | 0.000 | -1.139 | -0.671 |
| ===== | | | | | | |
| Omnibus: | 203.884 | | Durbin-Watson: | 0.921 | | |

| | | | |
|----------------|-------|-------------------|-----------|
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 1008.072 |
| Skew: | 1.728 | Prob(JB): | 1.26e-219 |
| Kurtosis: | 8.989 | Cond. No. | 1.34e+03 |

=====

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 1.34e+03. This might indicate that there are strong multicollinearity or other numerical problems.

0.0.20 Interpretation of multiple linear regression (lstat, rm, age, crim, ptratio)

The model explains about **68%** of the variation in housing prices ($R^2 = 0.683$), which is a big improvement over the single-variable models.

Intercept (const) → \$17.6k

This is the baseline price when all predictors are zero — mostly theoretical but useful for the model.

lstat (low socioeconomic households) → coef -0.58

As **lstat** increases by 1%, price drops about **\$580**. The negative effect stays strong even after controlling for other predictors.

rm (number of rooms) → coef $+4.46$

Each extra room adds about **\$4,460** to the price. Still a strong positive effect, though smaller than in the simple model.

age (percent old houses) → coef $+0.018$

This turns slightly positive but is **not significant** ($p = 0.082$). Age no longer matters much once we control for other factors.

crim (crime rate) → coef -0.07

Each unit increase in crime lowers price by about **\$70**. The effect is small but statistically significant ($p = 0.024$).

ptratio (pupil-teacher ratio) → coef -0.90

Each extra student per teacher reduces price by about **\$900**. The negative impact of crowded schools remains strong.

-
- The model improves overall fit **age** lost its importance when combined with others.

0.0.21 Fit medv as response with all available predictors altogether.

```
[20]: X = sm.add_constant(df.drop(columns='medv'))
model_full = sm.OLS(df['medv'], X).fit()
print(model_full.summary())
```

```

                                OLS Regression Results
=====
Dep. Variable:                  medv      R-squared:                0.741
Model:                          OLS      Adj. R-squared:           0.734
Method:                        Least Squares      F-statistic:           108.1
Date:                          Sat, 10 May 2025      Prob (F-statistic):       6.72e-135
Time:                           16:19:53      Log-Likelihood:          -1498.8
No. Observations:                506      AIC:                     3026.
Df Residuals:                    492      BIC:                     3085.
Df Model:                        13
Covariance Type:                  nonrobust
=====

```

| | coef | std err | t | P> t | [0.025 | 0.975] |
|---------|----------|---------|---------|-------|---------|---------|
| const | 36.4595 | 5.103 | 7.144 | 0.000 | 26.432 | 46.487 |
| crim | -0.1080 | 0.033 | -3.287 | 0.001 | -0.173 | -0.043 |
| zn | 0.0464 | 0.014 | 3.382 | 0.001 | 0.019 | 0.073 |
| indus | 0.0206 | 0.061 | 0.334 | 0.738 | -0.100 | 0.141 |
| chas | 2.6867 | 0.862 | 3.118 | 0.002 | 0.994 | 4.380 |
| nox | -17.7666 | 3.820 | -4.651 | 0.000 | -25.272 | -10.262 |
| rm | 3.8099 | 0.418 | 9.116 | 0.000 | 2.989 | 4.631 |
| age | 0.0007 | 0.013 | 0.052 | 0.958 | -0.025 | 0.027 |
| dis | -1.4756 | 0.199 | -7.398 | 0.000 | -1.867 | -1.084 |
| rad | 0.3060 | 0.066 | 4.613 | 0.000 | 0.176 | 0.436 |
| tax | -0.0123 | 0.004 | -3.280 | 0.001 | -0.020 | -0.005 |
| ptratio | -0.9527 | 0.131 | -7.283 | 0.000 | -1.210 | -0.696 |
| black | 0.0093 | 0.003 | 3.467 | 0.001 | 0.004 | 0.015 |
| lstat | -0.5248 | 0.051 | -10.347 | 0.000 | -0.624 | -0.425 |

```

=====
Omnibus:                        178.041      Durbin-Watson:              1.078
Prob(Omnibus):                  0.000      Jarque-Bera (JB):           783.126
Skew:                          1.521      Prob(JB):                   8.84e-171
Kurtosis:                      8.281      Cond. No.                   1.51e+04
=====

```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 1.51e+04. This might indicate that there are strong multicollinearity or other numerical problems.

0.0.22 Interpretation of full model with all predictors

This full model explains about **74%** of the variation in house prices ($R^2 = 0.741$), which is the best fit so far.

Intercept (const) → \$36.5k

Baseline price when all predictors are zero — mostly theoretical.

Key predictors and what they tell us:

- **crim** → -0.11 → Each unit increase in crime lowers price by ~\$110. Still a small but meaningful negative effect.
- **zn** → $+0.046$ → Higher residential zoning slightly raises prices (~\$46 per unit), small but significant.
- **indus** → not significant ($p = 0.738$) → industrial share has no clear price effect here.
- **chas** → $+2.69$ → Homes near the Charles River are ~\$2,690 higher on average.
- **nox** → -17.77 → Higher air pollution sharply lowers prices (~\$17.8k per unit).
- **rm** → $+3.81$ → Each extra room adds ~\$3,810, a strong positive effect.
- **age** → not significant ($p = 0.958$) → age loses influence in the full model.
- **dis** → -1.48 → Longer distance to employment centers cuts ~\$1,480 from price.
- **rad** → $+0.31$ → Better highway access slightly raises prices.
- **tax** → -0.012 → Higher taxes reduce price slightly (~\$12 per unit).
- **prratio** → -0.95 → More crowded schools lower price by ~\$950 per extra student per teacher.
- **black** → $+0.009$ → Higher Black population index slightly raises prices (~\$9 per unit), though the meaning is complex.
- **lstat** → -0.52 → Higher % low-status households still strongly reduces price (~\$520 per unit).

Notes: - Almost all predictors are significant ($p < 0.05$), except **indus** and **age**. - The **condition number is very high (~15,100)** → which tells us some predictors are tangled together, making it harder to trust the exact size of each effect.

0.0.23 Check the correlation between the predictors

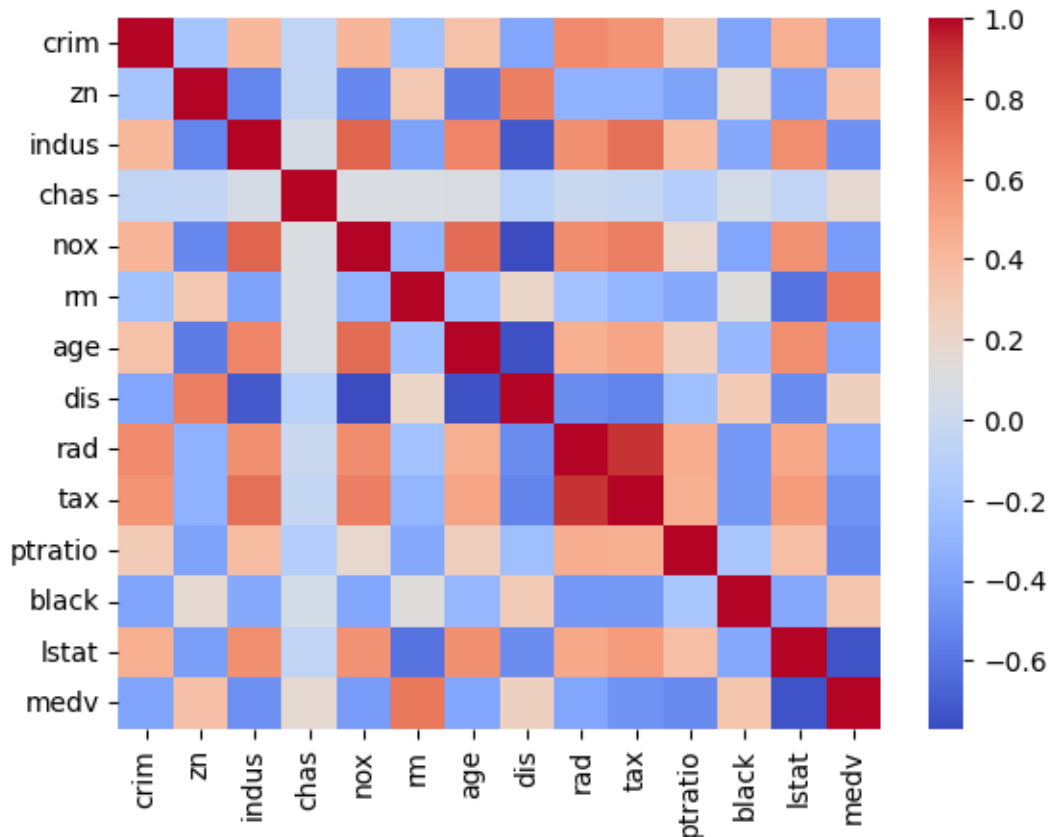
```
[21]: corr = df.corr()
      print(corr)

      # heat-map
      import seaborn as sns, matplotlib.pyplot as plt
      sns.heatmap(corr, cmap='coolwarm')
      plt.show()
```

| | crim | zn | indus | chas | nox | rm | age | \ |
|-------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|---|
| crim | 1.000000 | -0.200469 | 0.406583 | -0.055892 | 0.420972 | -0.219247 | 0.352734 | |
| zn | -0.200469 | 1.000000 | -0.533828 | -0.042697 | -0.516604 | 0.311991 | -0.569537 | |
| indus | 0.406583 | -0.533828 | 1.000000 | 0.062938 | 0.763651 | -0.391676 | 0.644779 | |

| | | | | | | | |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| chas | -0.055892 | -0.042697 | 0.062938 | 1.000000 | 0.091203 | 0.091251 | 0.086518 |
| nox | 0.420972 | -0.516604 | 0.763651 | 0.091203 | 1.000000 | -0.302188 | 0.731470 |
| rm | -0.219247 | 0.311991 | -0.391676 | 0.091251 | -0.302188 | 1.000000 | -0.240265 |
| age | 0.352734 | -0.569537 | 0.644779 | 0.086518 | 0.731470 | -0.240265 | 1.000000 |
| dis | -0.379670 | 0.664408 | -0.708027 | -0.099176 | -0.769230 | 0.205246 | -0.747881 |
| rad | 0.625505 | -0.311948 | 0.595129 | -0.007368 | 0.611441 | -0.209847 | 0.456022 |
| tax | 0.582764 | -0.314563 | 0.720760 | -0.035587 | 0.668023 | -0.292048 | 0.506456 |
| ptratio | 0.289946 | -0.391679 | 0.383248 | -0.121515 | 0.188933 | -0.355501 | 0.261515 |
| black | -0.385064 | 0.175520 | -0.356977 | 0.048788 | -0.380051 | 0.128069 | -0.273534 |
| lstat | 0.455621 | -0.412995 | 0.603800 | -0.053929 | 0.590879 | -0.613808 | 0.602339 |
| medv | -0.388305 | 0.360445 | -0.483725 | 0.175260 | -0.427321 | 0.695360 | -0.376955 |

| | | | | | | | |
|---------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| | dis | rad | tax | ptratio | black | lstat | medv |
| crim | -0.379670 | 0.625505 | 0.582764 | 0.289946 | -0.385064 | 0.455621 | -0.388305 |
| zn | 0.664408 | -0.311948 | -0.314563 | -0.391679 | 0.175520 | -0.412995 | 0.360445 |
| indus | -0.708027 | 0.595129 | 0.720760 | 0.383248 | -0.356977 | 0.603800 | -0.483725 |
| chas | -0.099176 | -0.007368 | -0.035587 | -0.121515 | 0.048788 | -0.053929 | 0.175260 |
| nox | -0.769230 | 0.611441 | 0.668023 | 0.188933 | -0.380051 | 0.590879 | -0.427321 |
| rm | 0.205246 | -0.209847 | -0.292048 | -0.355501 | 0.128069 | -0.613808 | 0.695360 |
| age | -0.747881 | 0.456022 | 0.506456 | 0.261515 | -0.273534 | 0.602339 | -0.376955 |
| dis | 1.000000 | -0.494588 | -0.534432 | -0.232471 | 0.291512 | -0.496996 | 0.249929 |
| rad | -0.494588 | 1.000000 | 0.910228 | 0.464741 | -0.444413 | 0.488676 | -0.381626 |
| tax | -0.534432 | 0.910228 | 1.000000 | 0.460853 | -0.441808 | 0.543993 | -0.468536 |
| ptratio | -0.232471 | 0.464741 | 0.460853 | 1.000000 | -0.177383 | 0.374044 | -0.507787 |
| black | 0.291512 | -0.444413 | -0.441808 | -0.177383 | 1.000000 | -0.366087 | 0.333461 |
| lstat | -0.496996 | 0.488676 | 0.543993 | 0.374044 | -0.366087 | 1.000000 | -0.737663 |
| medv | 0.249929 | -0.381626 | -0.468536 | -0.507787 | 0.333461 | -0.737663 | 1.000000 |



0.0.24 Interpretation of the correlation matrix

Looking at the correlation matrix and heatmap, we can see how strongly the predictors are related to each other and to the response `medv`.

Strong correlations with house price (`medv`): - `lstat` $\rightarrow -0.74$ \rightarrow strong negative \rightarrow more low-status households = lower price. - `rm` $\rightarrow +0.70$ \rightarrow strong positive \rightarrow more rooms = higher price. - `ptratio` $\rightarrow -0.51$ \rightarrow moderately negative \rightarrow larger class sizes = lower price. - `indus`, `nox`, `crim` \rightarrow moderate negative \rightarrow more industry, pollution, or crime = lower price. - `black`, `zn` \rightarrow weak positive.

Strong correlations between predictors (possible multicollinearity): - `rad` and `tax` $\rightarrow +0.91$ \rightarrow more highway access = higher taxes. - `nox` and `indus` $\rightarrow +0.76$ \rightarrow more industry = more pollution. - `dis` and `nox` $\rightarrow -0.77$ \rightarrow greater distance from jobs = less pollution. - `age` and `dis` $\rightarrow -0.75$ \rightarrow older areas are closer to city centers.

What this tells us: - Some predictors are tightly linked (like `rad` and `tax`), so in regression models

we need to **watch out for multicollinearity**. - The strongest price drivers (`lstat` and `rm`) stand out clearly. - Pollution, industry, crime, and school quality all matter, but they're also tangled with each other.

0.0.25 Use the multiple linear regression model

```
[25]: import itertools

# Define the values
lstatC = [5, 10, 15]
rmC = [5, 6.5, 8]

# Create all combinations (expand.grid equivalent)
selected_predictor_values = pd.DataFrame(list(itertools.product(lstatC, rmC)),
                                         columns=['lstat', 'rm'])

# Show the grid
print(selected_predictor_values)
```

| | lstat | rm |
|---|-------|-----|
| 0 | 5 | 5.0 |
| 1 | 5 | 6.5 |
| 2 | 5 | 8.0 |
| 3 | 10 | 5.0 |
| 4 | 10 | 6.5 |
| 5 | 10 | 8.0 |
| 6 | 15 | 5.0 |
| 7 | 15 | 6.5 |
| 8 | 15 | 8.0 |

```
[26]: import statsmodels.api as sm

# Fit the model on lstat + rm
X = sm.add_constant(df[['lstat', 'rm']])
model = sm.OLS(df['medv'], X).fit()

# Add constant to prediction grid
grid_with_const = sm.add_constant(selected_predictor_values)

# Predict with intervals
predictions = model.get_prediction(grid_with_const)
print(predictions.summary_frame(alpha=0.05))
```

| | mean | mean_se | mean_ci_lower | mean_ci_upper | obs_ci_lower | \ |
|---|-----------|----------|---------------|---------------|--------------|---|
| 0 | 20.903875 | 0.856315 | 19.221481 | 22.586269 | 9.889729 | |
| 1 | 28.546057 | 0.377499 | 27.804387 | 29.287727 | 17.635923 | |
| 2 | 36.188239 | 0.663860 | 34.883959 | 37.492519 | 25.225479 | |
| 3 | 17.692084 | 0.693873 | 16.328837 | 19.055330 | 6.722152 | |

| | | | | | |
|---|-----------|----------|-----------|-----------|-----------|
| 4 | 25.334266 | 0.263915 | 24.815754 | 25.852777 | 14.437027 |
| 5 | 32.976448 | 0.739470 | 31.523618 | 34.429277 | 21.995024 |
| 6 | 14.480292 | 0.570322 | 13.359785 | 15.600799 | 3.537875 |
| 7 | 22.122474 | 0.304004 | 21.525200 | 22.719748 | 11.221204 |
| 8 | 29.764656 | 0.865184 | 28.064837 | 31.464475 | 18.747835 |

| | obs_ci_upper |
|---|--------------|
| 0 | 31.918021 |
| 1 | 39.456192 |
| 2 | 47.150999 |
| 3 | 28.662016 |
| 4 | 36.231505 |
| 5 | 43.957872 |
| 6 | 25.422709 |
| 7 | 33.023745 |
| 8 | 40.781477 |

0.0.26 Interpretation of predictions for `lstat` and `rm`

We predicted median house prices (`medv`) for combinations of `lstat` (% low-status households) and `rm` (average number of rooms), along with 95% prediction intervals.

Summary of patterns:

- **Low `lstat`, high `rm` → highest prices:**
Example: `lstat` = 5, `rm` = 8 → predicted ~\$36.2k, range ~\$25.2k–\$47.2k.
- **High `lstat`, low `rm` → lowest prices:**
Example: `lstat` = 15, `rm` = 5 → predicted ~\$14.5k, range ~\$3.5k–\$25.4k.
- **Effect of `lstat` at fixed `rm`:**
Prices drop as `lstat` increases.
Example at `rm` = 6.5 →
`lstat` = 5 → ~\$28.5k,
`lstat` = 10 → ~\$25.3k,
`lstat` = 15 → ~\$22.1k.
- **Effect of `rm` at fixed `lstat`:**
Prices rise as rooms increase.
Example at `lstat` = 10 →
`rm` = 5 → ~\$17.7k,
`rm` = 6.5 → ~\$25.3k,
`rm` = 8 → ~\$33.0k.

About the prediction intervals: - Across all combinations, the prediction intervals are fairly similar in width (~\$22k). - This means the model's uncertainty stays consistent whether predicting low-end or high-end prices.

in simple terms:

Low-status, small homes \rightarrow lowest predicted prices.

Low-status, large homes \rightarrow middle range.

High-status, large homes \rightarrow highest prices.

Prediction intervals help us see the expected spread in prices for similar homes.