

# Assignments 2+3 - Simple and Multiple Linear Regression (I+II)

## A. Conceptual Questions

Suppose we have a data set with five predictors,  $X_1 = \text{GPA}$ ,  $X_2 = \text{IQ}$ ,  $X_3 = \text{Level}$  (1 for College and 0 for High School),  $X_4 = \text{Interaction between GPA and IQ}$ , and  $X_5 = \text{Interaction between GPA and Level}$ . The response is starting salary after graduation (in thousands of dollars). Suppose we use least squares to fit the model, and get  $\hat{\beta}_0 = 50$ ,  $\hat{\beta}_1 = 20$ ,  $\hat{\beta}_2 = 0.07$ ,  $\hat{\beta}_3 = 35$ ,  $\hat{\beta}_4 = 0.01$ ,  $\hat{\beta}_5 = -10$ .

### 1. Which answer is correct, and why?

- i. For a fixed value of IQ and GPA, high school graduates earn more, on average, than college graduates.
- ii. For a fixed value of IQ and GPA, college graduates earn more, on average, than high school graduates.
- iii. For a fixed value of IQ and GPA, high school graduates earn more, on average, than college graduates provided that the GPA is high enough.
- iv. For a fixed value of IQ and GPA, college graduates earn more, on average, than high school graduates provided that the GPA is high enough.

---

--- Your answer here ---

**Hint:** Use latex code to write your equations.

---

### 2. Predict the salary of a college graduate with IQ of 110 and a GPA of 4.0.

---

--- Your answer here ---

**Hint:** Use latex code to write your equations.

---

### 3. True or false: Since the coefficient for the GPA/IQ interaction term is very small, there is very little evidence of an interaction effect. Justify your answer.

---

--- Your answer here ---

## B. Practical

### Overview of the steps

Assignment 2:

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

Assignment 3: 6. Add interaction terms 7. Apply non-linear transformations to some predictors 8. Use categorical predictors

### Steps of Assignment 2 in detail

#### Load the data and get an overview of the data

Load the data file `Boston.rda` or `Boston.csv` .

In R the dataframe comes with the `MASS` library. We save the dataframe ones in `csv` and `rda` files for later use.

```
In [74]: 1 library(MASS)
          2 #write.csv(Boston,"../ISLR/data/Boston.csv", row.names = TRUE)
          3 #save(Boston,file="../ISLR/data/Boston.rda")
```

Display the number of predictors (including the response `medv` ) and their names:

```
In [75]: 1 dim(Boston)[2]
          2 names(Boston)
```

14

'crim' 'zn' 'indus' 'chas' 'nox' 'rm' 'age' 'dis' 'rad' 'tax' 'ptratio' 'black'  
'lstat' 'medv'

Print a statistic summary of the predictors and the response `medv` :

In [76]:

1 summary(Boston)

crim		zn		indus		chas	
Min.	: 0.00632	Min.	: 0.00	Min.	: 0.46	Min.	:0.00000
1st Qu.:	0.08204	1st Qu.:	0.00	1st Qu.:	5.19	1st Qu.:	0.00000
Median	: 0.25651	Median	: 0.00	Median	: 9.69	Median	:0.00000
Mean	: 3.61352	Mean	: 11.36	Mean	:11.14	Mean	:0.06917
3rd Qu.:	3.67708	3rd Qu.:	12.50	3rd Qu.:	18.10	3rd Qu.:	0.00000
Max.	:88.97620	Max.	:100.00	Max.	:27.74	Max.	:1.00000
nox		rm		age		dis	
Min.	:0.3850	Min.	:3.561	Min.	: 2.90	Min.	: 1.130
1st Qu.:	0.4490	1st Qu.:	5.886	1st Qu.:	45.02	1st Qu.:	2.100
Median	:0.5380	Median	:6.208	Median	: 77.50	Median	: 3.207
Mean	:0.5547	Mean	:6.285	Mean	: 68.57	Mean	: 3.795
3rd Qu.:	0.6240	3rd Qu.:	6.623	3rd Qu.:	94.08	3rd Qu.:	5.188
Max.	:0.8710	Max.	:8.780	Max.	:100.00	Max.	:12.127
rad		tax		ptratio		black	
Min.	: 1.000	Min.	:187.0	Min.	:12.60	Min.	: 0.32
1st Qu.:	4.000	1st Qu.:	279.0	1st Qu.:	17.40	1st Qu.:	375.38
Median	: 5.000	Median	:330.0	Median	:19.05	Median	:391.44
Mean	: 9.549	Mean	:408.2	Mean	:18.46	Mean	:356.67
3rd Qu.:	24.000	3rd Qu.:	666.0	3rd Qu.:	20.20	3rd Qu.:	396.23
Max.	:24.000	Max.	:711.0	Max.	:22.00	Max.	:396.90
lstat		medv					
Min.	: 1.73	Min.	: 5.00				
1st Qu.:	6.95	1st Qu.:	17.02				
Median	:11.36	Median	:21.20				
Mean	:12.65	Mean	:22.53				
3rd Qu.:	16.95	3rd Qu.:	25.00				
Max.	:37.97	Max.	:50.00				

Display the number of data points:

In [77]:

1 dim(Boston)[1]

506

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

In [78]:

1	Boston
---	--------

A data.frame: 506 × 14

	crim	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	black
	<dbl>	<dbl>	<dbl>	<int>	<dbl>	<dbl>	<dbl>	<dbl>	<int>	<dbl>	<dbl>	<dbl>
1	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296	15.3	396.90
2	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.90
3	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242	17.8	392.83
4	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63
5	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	396.90
6	0.02985	0.0	2.18	0	0.458	6.430	58.7	6.0622	3	222	18.7	394.12
7	0.08829	12.5	7.87	0	0.524	6.012	66.6	5.5605	5	311	15.2	395.60
8	0.14455	12.5	7.87	0	0.524	6.172	96.1	5.9505	5	311	15.2	396.90
9	0.21124	12.5	7.87	0	0.524	5.631	100.0	6.0821	5	311	15.2	386.63
10	0.17004	12.5	7.87	0	0.524	6.004	85.9	6.5921	5	311	15.2	386.71
11	0.22489	12.5	7.87	0	0.524	6.377	94.3	6.3467	5	311	15.2	392.52
12	0.11747	12.5	7.87	0	0.524	6.009	82.9	6.2267	5	311	15.2	396.90
13	0.09378	12.5	7.87	0	0.524	5.889	39.0	5.4509	5	311	15.2	390.50
14	0.62976	0.0	8.14	0	0.538	5.949	61.8	4.7075	4	307	21.0	396.90
15	0.63796	0.0	8.14	0	0.538	6.096	84.5	4.4619	4	307	21.0	380.02
16	0.62739	0.0	8.14	0	0.538	5.834	56.5	4.4986	4	307	21.0	395.62
17	1.05393	0.0	8.14	0	0.538	5.935	29.3	4.4986	4	307	21.0	386.85
18	0.78420	0.0	8.14	0	0.538	5.990	81.7	4.2579	4	307	21.0	386.75
19	0.80271	0.0	8.14	0	0.538	5.456	36.6	3.7965	4	307	21.0	288.99
20	0.72580	0.0	8.14	0	0.538	5.727	69.5	3.7965	4	307	21.0	390.95
21	1.25179	0.0	8.14	0	0.538	5.570	98.1	3.7979	4	307	21.0	376.57
22	0.85204	0.0	8.14	0	0.538	5.965	89.2	4.0123	4	307	21.0	392.53
23	1.23247	0.0	8.14	0	0.538	6.142	91.7	3.9769	4	307	21.0	396.90
24	0.98843	0.0	8.14	0	0.538	5.813	100.0	4.0952	4	307	21.0	394.54
25	0.75026	0.0	8.14	0	0.538	5.924	94.1	4.3996	4	307	21.0	394.33
26	0.84054	0.0	8.14	0	0.538	5.599	85.7	4.4546	4	307	21.0	303.42
27	0.67191	0.0	8.14	0	0.538	5.813	90.3	4.6820	4	307	21.0	376.88
28	0.95577	0.0	8.14	0	0.538	6.047	88.8	4.4534	4	307	21.0	306.38
29	0.77299	0.0	8.14	0	0.538	6.495	94.4	4.4547	4	307	21.0	387.94
30	1.00245	0.0	8.14	0	0.538	6.674	87.3	4.2390	4	307	21.0	380.23
:	:	:	:	:	:	:	:	:	:	:	:	:
477	4.87141	0	18.10	0	0.614	6.484	93.6	2.3053	24	666	20.2	396.21
478	15.02340	0	18.10	0	0.614	5.304	97.3	2.1007	24	666	20.2	349.48
479	10.23300	0	18.10	0	0.614	6.185	96.7	2.1705	24	666	20.2	379.70
480	14.33370	0	18.10	0	0.614	6.229	88.0	1.9512	24	666	20.2	383.32

	crim	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	black
	<dbl>	<dbl>	<dbl>	<int>	<dbl>	<dbl>	<dbl>	<dbl>	<int>	<dbl>	<dbl>	<dbl>
481	5.82401	0	18.10	0	0.532	6.242	64.7	3.4242	24	666	20.2	396.90
482	5.70818	0	18.10	0	0.532	6.750	74.9	3.3317	24	666	20.2	393.07
483	5.73116	0	18.10	0	0.532	7.061	77.0	3.4106	24	666	20.2	395.28
484	2.81838	0	18.10	0	0.532	5.762	40.3	4.0983	24	666	20.2	392.92
485	2.37857	0	18.10	0	0.583	5.871	41.9	3.7240	24	666	20.2	370.73
486	3.67367	0	18.10	0	0.583	6.312	51.9	3.9917	24	666	20.2	388.62
487	5.69175	0	18.10	0	0.583	6.114	79.8	3.5459	24	666	20.2	392.68
488	4.83567	0	18.10	0	0.583	5.905	53.2	3.1523	24	666	20.2	388.22
489	0.15086	0	27.74	0	0.609	5.454	92.7	1.8209	4	711	20.1	395.09
490	0.18337	0	27.74	0	0.609	5.414	98.3	1.7554	4	711	20.1	344.05
491	0.20746	0	27.74	0	0.609	5.093	98.0	1.8226	4	711	20.1	318.43
492	0.10574	0	27.74	0	0.609	5.983	98.8	1.8681	4	711	20.1	390.11
493	0.11132	0	27.74	0	0.609	5.983	83.5	2.1099	4	711	20.1	396.90
494	0.17331	0	9.69	0	0.585	5.707	54.0	2.3817	6	391	19.2	396.90
495	0.27957	0	9.69	0	0.585	5.926	42.6	2.3817	6	391	19.2	396.90
496	0.17899	0	9.69	0	0.585	5.670	28.8	2.7986	6	391	19.2	393.29
497	0.28960	0	9.69	0	0.585	5.390	72.9	2.7986	6	391	19.2	396.90
498	0.26838	0	9.69	0	0.585	5.794	70.6	2.8927	6	391	19.2	396.90
499	0.23912	0	9.69	0	0.585	6.019	65.3	2.4091	6	391	19.2	396.90
500	0.17783	0	9.69	0	0.585	5.569	73.5	2.3999	6	391	19.2	395.77
501	0.22438	0	9.69	0	0.585	6.027	79.7	2.4982	6	391	19.2	396.90
502	0.06263	0	11.93	0	0.573	6.593	69.1	2.4786	1	273	21.0	391.99
503	0.04527	0	11.93	0	0.573	6.120	76.7	2.2875	1	273	21.0	396.90
504	0.06076	0	11.93	0	0.573	6.976	91.0	2.1675	1	273	21.0	396.90
505	0.10959	0	11.93	0	0.573	6.794	89.3	2.3889	1	273	21.0	393.45
506	0.04741	0	11.93	0	0.573	6.030	80.8	2.5050	1	273	21.0	396.90

Plot some predictors (at least two) against the response values. We choose `lstat` , `rm` , and `age` .

In `R` , we need to download and install a library first.

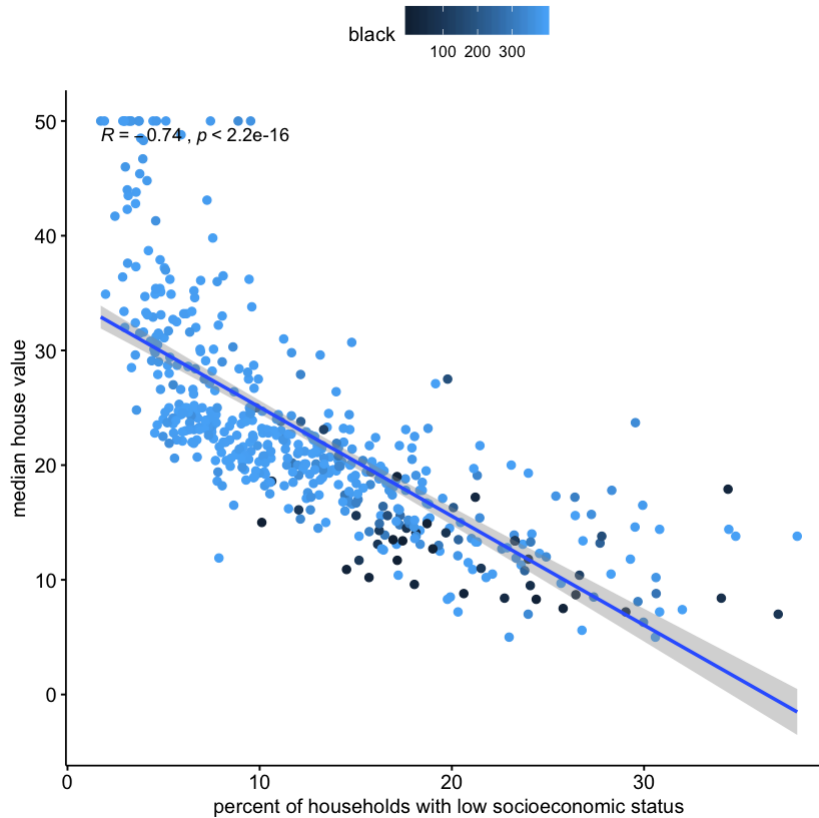
In [79]:

```
1 install.packages("ggpubr")
2 library("ggpubr")
```

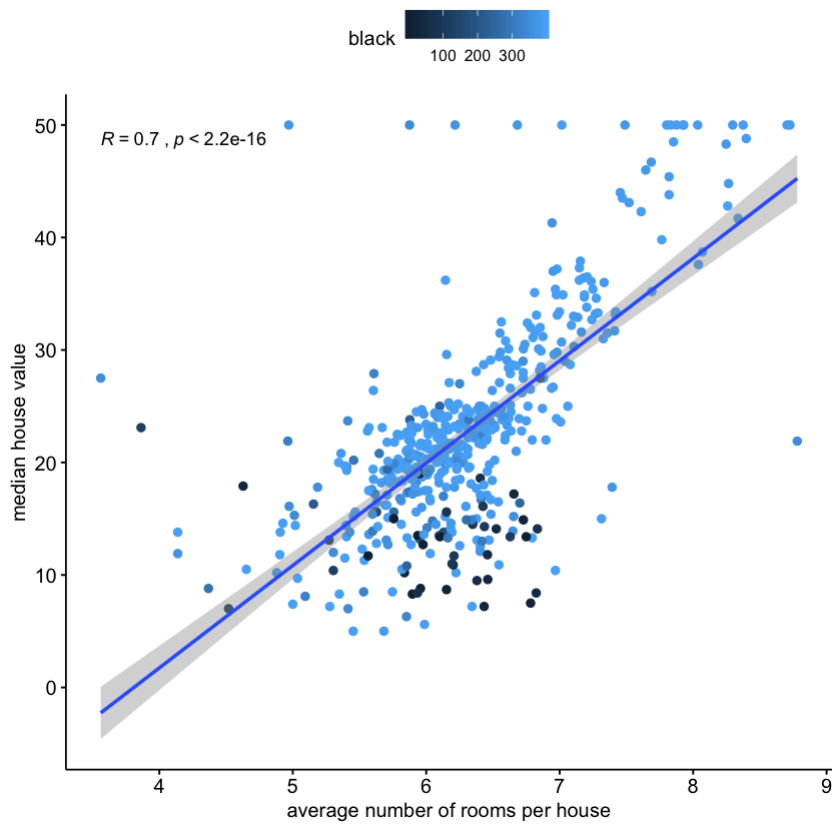
The downloaded binary packages are in  
/var/folders/ct/4pcck8t94sdfc73rhymq4t140000gp/T//RtmpiV0p6  
t/downloaded\_packages

The R function `ggscatter` even displays a regression line, confidence intervals, the Pearson coefficient of correlation, and the  $p$  value. **This is not necessary at this stage.**

```
In [80]: 1 ggscatter(Boston, x = "lstat", y = "medv",  
2         add = "reg.line", conf.int = TRUE,  
3         cor.coef = TRUE, cor.method = "pearson",  
4         xlab = "percent of households with low socioeconomic sta  
5         ylab = "median house value")
```

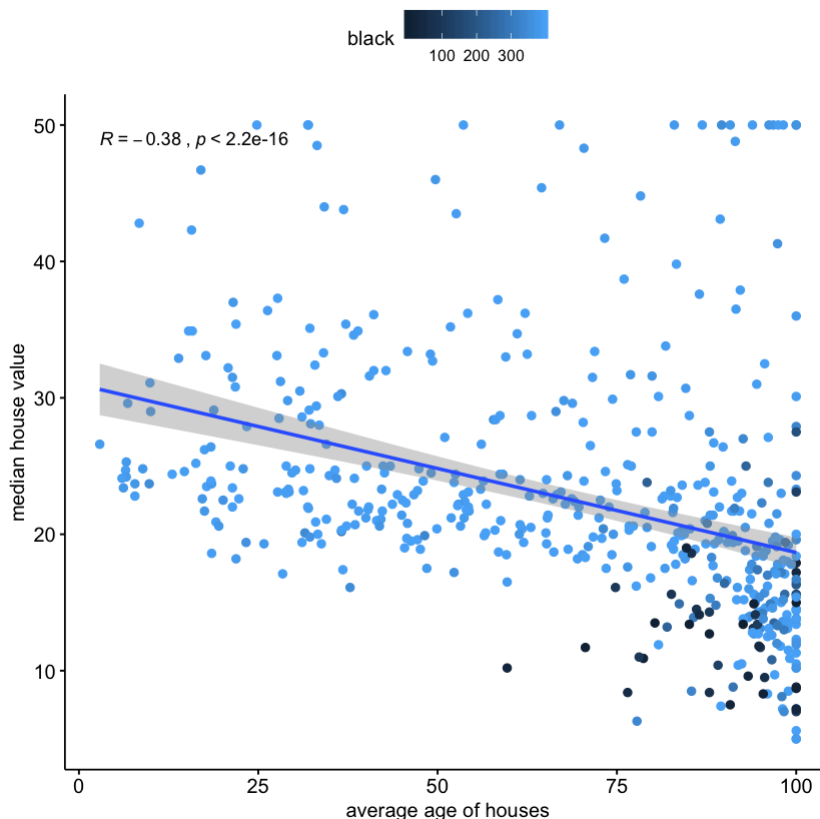


```
In [81]: 1 ggscatter(Boston, x = "rm", y = "medv",  
2         add = "reg.line", conf.int = TRUE,  
3         cor.coef = TRUE, cor.method = "pearson",  
4         xlab = "average number of rooms per house",  
5         ylab = "median house value")
```





```
In [82]: 1 ggscatter(Boston, x = "age", y = "medv",
2           add = "reg.line", conf.int = TRUE,
3           cor.coef = TRUE, cor.method = "pearson",
4           xlab = "average age of houses",
5           ylab = "median house value")
```



## Perform simple linear regressions

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

```
In [83]: 1 lm.fit=lm(medv~lstat ,data=Boston)
2 summary(lm.fit)
```

Call:

```
lm(formula = medv ~ lstat, data = Boston)
```

Residuals:

Min	1Q	Median	3Q	Max
-15.168	-3.990	-1.318	2.034	24.500

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	34.55384	0.56263	61.41	<2e-16 ***
lstat	-0.95005	0.03873	-24.53	<2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.216 on 504 degrees of freedom  
 Multiple R-squared: 0.5441, Adjusted R-squared: 0.5432  
 F-statistic: 601.6 on 1 and 504 DF, p-value: < 2.2e-16

In [84]:

```
1 lm.fit=lm(medv~rm ,data=Boston)
2 summary(lm.fit)
```

Call:  
lm(formula = medv ~ rm, data = Boston)

Residuals:

Min	1Q	Median	3Q	Max
-23.346	-2.547	0.090	2.986	39.433

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-34.671	2.650	-13.08	<2e-16 ***
rm	9.102	0.419	21.72	<2e-16 ***

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.616 on 504 degrees of freedom  
Multiple R-squared: 0.4835, Adjusted R-squared: 0.4825  
F-statistic: 471.8 on 1 and 504 DF, p-value: < 2.2e-16

In [85]:

```
1 lm.fit=lm(medv~age ,data=Boston)
2 summary(lm.fit)
```

Call:  
lm(formula = medv ~ age, data = Boston)

Residuals:

Min	1Q	Median	3Q	Max
-15.097	-5.138	-1.958	2.397	31.338

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	30.97868	0.99911	31.006	<2e-16 ***
age	-0.12316	0.01348	-9.137	<2e-16 ***

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8.527 on 504 degrees of freedom  
Multiple R-squared: 0.1421, Adjusted R-squared: 0.1404  
F-statistic: 83.48 on 1 and 504 DF, p-value: < 2.2e-16

Interprete the results. *Your interpretation of the results goes here!*

Obtain a confidence interval for the coefficient estimates for the indivisual models.

In [86]:

```
1 lm.fit=lm(medv~lstat ,data=Boston)
2 confint(lm.fit)
```

A matrix: 2 × 2 of type dbl

	2.5 %	97.5 %
(Intercept)	33.448457	35.6592247
lstat	-1.026148	-0.8739505

```
In [87]: 1 lm.fit=lm(medv~rm ,data=Boston)
2 confint(lm.fit)
```

A matrix: 2 × 2 of type dbl

	2.5 %	97.5 %
(Intercept)	-39.876641	-29.464601
rm	8.278855	9.925363

```
In [88]: 1 lm.fit=lm(medv~age ,data=Boston)
2 confint(lm.fit)
```

A matrix: 2 × 2 of type dbl

	2.5 %	97.5 %
(Intercept)	29.0157516	32.94160395
age	-0.1496469	-0.09667852

Interprete the results. *Your interpretation of the results goes here!*

### Use the simple linear regression models

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

```
In [89]: 1 lm.fit=lm(medv~lstat,data=Boston)
2 predict(lm.fit,data.frame(lstat=c(5,10,15)), interval ="prediction")
```

A matrix: 3 × 3 of type dbl

	fit	lwr	upr
1	29.80359	17.565675	42.04151
2	25.05335	12.827626	37.27907
3	20.30310	8.077742	32.52846

```
In [90]: 1 lm.fit=lm(medv~rm,data=Boston)
2 predict(lm.fit,data.frame(rm=c(5,6.5,8)), interval ="prediction")
```

A matrix: 3 × 3 of type dbl

	fit	lwr	upr
1	10.83992	-2.214474	23.89432
2	24.49309	11.480391	37.50578
3	38.14625	25.058353	51.23415

In [91]:

1

2

lm.fit=lm(medv~age,data=Boston)  
predict(lm.fit,data.frame(age=c(25,50,75)), interval ="prediction'

A matrix: 3 × 3 of type dbl

	fit	lwr	upr
1	27.89961	11.090368	44.70885
2	24.82054	8.043748	41.59734
3	21.74147	4.971031	38.51192

Interprete the results. *Your interpretation of the results goes here!*

Perform multiple linear regressions

Fit medv as response with the predictors selected before altogether.

In [92]:

1

2

lm.fit=lm(medv~lstat+rm+age ,data=Boston)  
summary(lm.fit)

Call:  
lm(formula = medv ~ lstat + rm + age, data = Boston)

Residuals:

Min	1Q	Median	3Q	Max
-18.210	-3.467	-1.053	1.957	27.500

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-1.175311	3.181924	-0.369	0.712
lstat	-0.668513	0.054357	-12.298	<2e-16 ***
rm	5.019133	0.454306	11.048	<2e-16 ***
age	0.009091	0.011215	0.811	0.418

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 5.542 on 502 degrees of freedom  
Multiple R-squared: 0.639, Adjusted R-squared: 0.6369  
F-statistic: 296.2 on 3 and 502 DF, p-value: < 2.2e-16

Interprete the results. *Your interpretation of the results goes here!*

Fit medv as response with all available predictors altogether.

```
In [93]: 1 lm.fit=lm(medv~. ,data=Boston)
          2 summary(lm.fit)
```

Call:

```
lm(formula = medv ~ ., data = Boston)
```

Residuals:

```
      Min       1Q   Median       3Q      Max
-15.595  -2.730  -0.518   1.777   26.199
```

Coefficients:

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  3.646e+01  5.103e+00   7.144 3.28e-12 ***
crim         -1.080e-01  3.286e-02  -3.287 0.001087 **
zn           4.642e-02  1.373e-02   3.382 0.000778 ***
indus        2.056e-02  6.150e-02   0.334 0.738288
chas         2.687e+00  8.616e-01   3.118 0.001925 **
nox          -1.777e+01  3.820e+00  -4.651 4.25e-06 ***
rm           3.810e+00  4.179e-01   9.116 < 2e-16 ***
age          6.922e-04  1.321e-02   0.052 0.958229
dis          -1.476e+00  1.995e-01  -7.398 6.01e-13 ***
rad           3.060e-01  6.635e-02   4.613 5.07e-06 ***
tax          -1.233e-02  3.760e-03  -3.280 0.001112 **
ptratio      -9.527e-01  1.308e-01  -7.283 1.31e-12 ***
black         9.312e-03  2.686e-03   3.467 0.000573 ***
lstat        -5.248e-01  5.072e-02 -10.347 < 2e-16 ***
---

```

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.745 on 492 degrees of freedom

Multiple R-squared: 0.7406, Adjusted R-squared: 0.7338

F-statistic: 108.1 on 13 and 492 DF, p-value: < 2.2e-16

Interprete the results. *Your interpretation of the results goes here!*

```
In [94]: 1 install.packages("corrplot")
          2 source("http://www.sthda.com/upload/rquery_cormat.r")
```

The downloaded binary packages are in

```
/var/folders/ct/4pcck8t94sdfc73rhymq4t140000gp/T//RtmpiV0p6
t/downloaded_packages
```

Check the correlation between the predictors.

In R , we need to download and install a library and an external function first.

In [95]:

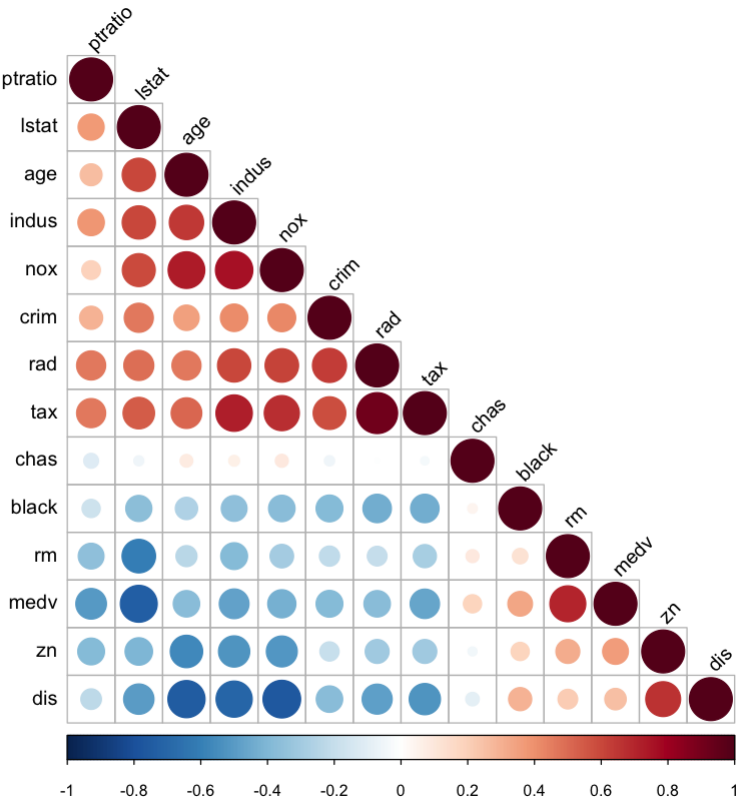
1	<code>rquery.cormat(Boston)</code>
---	------------------------------------

\$r		ptratio	lstat	age	indus	nox	crim	rad	tax	chas
s	black									
ptratio		1								
lstat		0.37	1							
age		0.26	0.6	1						
indus		0.38	0.6	0.64	1					
nox		0.19	0.59	0.73	0.76	1				
crim		0.29	0.46	0.35	0.41	0.42	1			
rad		0.46	0.49	0.46	0.6	0.61	0.63	1		
tax		0.46	0.54	0.51	0.72	0.67	0.58	0.91	1	
chas		-0.12	-0.054	0.087	0.063	0.091	-0.056	-0.0074	-0.036	
1	black	-0.18	-0.37	-0.27	-0.36	-0.38	-0.39	-0.44	-0.44	0.04
9	1	-0.36	-0.61	-0.24	-0.39	-0.3	-0.22	-0.21	-0.29	0.09
rm		-0.36	-0.61	-0.24	-0.39	-0.3	-0.22	-0.21	-0.29	0.09
1	0.13	-0.51	-0.74	-0.38	-0.48	-0.43	-0.39	-0.38	-0.47	0.1
medv		-0.51	-0.74	-0.38	-0.48	-0.43	-0.39	-0.38	-0.47	0.1
8	0.33	-0.39	-0.41	-0.57	-0.53	-0.52	-0.2	-0.31	-0.31	-0.04
zn		-0.39	-0.41	-0.57	-0.53	-0.52	-0.2	-0.31	-0.31	-0.04
3	0.18	-0.23	-0.5	-0.75	-0.71	-0.77	-0.38	-0.49	-0.53	-0.09
dis		-0.23	-0.5	-0.75	-0.71	-0.77	-0.38	-0.49	-0.53	-0.09
9	0.29									
	rm medv			zn	dis					
ptratio										
lstat										
age										
indus										
nox										
crim										
rad										
tax										
chas										
black										
rm		1								
medv		0.7	1							
zn		0.31	0.36	1						
dis		0.21	0.25	0.66	1					
\$p										
	ptratio		lstat	age	indus	nox	crim	rad		
tax										
ptratio		0								
lstat		3e-18	0							
age		2.3e-09	2.8e-51	0						
indus		3.8e-19	1.4e-51	8.4e-61	0					
nox		1.9e-05	6e-49	7.5e-86	7.9e-98	0				
crim		2.9e-11	2.7e-27	2.9e-16	1.5e-21	3.8e-23	0			
rad		1.8e-28	9.9e-32	2.4e-27	8.4e-50	3.3e-53	2.7e-56	0		
tax		5.7e-28	2.6e-40	2.6e-34	3e-82	1.1e-66	2.4e-47	4.1e-195		
0	chas	0.0062	0.23	0.052	0.16	0.04	0.21	0.87		
0.42	black	6e-05	1.7e-17	3.9e-10	1.2e-16	7.8e-19	2.5e-19	6.6e-26	1.	
4e-25	rm	1.6e-16	1e-53	4.5e-08	5.3e-20	3.8e-12	6.3e-07	1.9e-06	2.	
1e-11	medv	1.6e-34	5.1e-88	1.6e-18	4.9e-31	7.1e-24	1.2e-19	5.5e-19	5.	
6e-29	zn	5.3e-20	2.9e-22	7.6e-45	1.3e-38	7.2e-36	5.5e-06	7e-13	4.	
4e-13										

```
dis      1.2e-07 6.4e-33 9.9e-92 3.6e-78 4.2e-100 8.5e-19 1.4e-32
1e-38
      chas    black      rm    medv      zn dis
ptratio
lstat
age
indus
nox
crim
rad
tax
chas      0
black    0.27      0
rm      0.04 0.0039      0
medv    7.4e-05 1.3e-14 2.5e-74      0
zn      0.34 7.2e-05 6.9e-13 5.7e-17      0
dis      0.026 2.3e-11 3.2e-06 1.2e-08 9.7e-66 0

$sym
      ptratio lstat age indus nox crim rad tax chas black rm medv
zn dis
ptratio 1
lstat .      1
age      .      1
indus .      .      ,      1
nox      .      .      ,      ,      1
crim     .      .      .      .      1
rad      .      .      .      .      ,      ,      1
tax      .      .      .      ,      ,      .      *      1
chas     .      .      .      .      .      .      .      1
black    .      .      .      .      .      .      .      1
rm      .      ,      .      .      .      .      .      .      1
medv     .      ,      .      .      .      .      .      .      ,      1
zn      .      .      .      .      .      .      .      .      .      .
1
dis      .      .      ,      ,      ,      .      .      .
, 1
attr("legend")
[1] 0 ' ' 0.3 '.' 0.6 ',' 0.8 '+' 0.9 '*' 0.95 'B' 1
```





Interprete the results. *Your interpretation of the results goes here!*

Use the multiple linear regression model

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

```
In [96]: 1 lstatC=c(5,10,15)
2 rmC=c(5,6.5,8)
3 selected_predictor_values = expand.grid(lstat = lstatC, rm = rmC)
4 selected_predictor_values
```

A data.frame: 9  
× 2

lstat	rm
<dbl>	<dbl>
5	5.0
10	5.0
15	5.0
5	6.5
10	6.5
15	6.5
5	8.0
10	8.0
15	8.0

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

```
In [97]: 1 lm.fit=lm(medv~lstat+rm ,data=Boston)
2 predict(lm.fit, selected_predictor_values, interval ="prediction")
```

A matrix: 9 × 3 of type dbl

	fit	lwr	upr
1	20.90388	9.889729	31.91802
2	17.69208	6.722152	28.66202
3	14.48029	3.537875	25.42271
4	28.54606	17.635923	39.45619
5	25.33427	14.437027	36.23150
6	22.12247	11.221204	33.02374
7	36.18824	25.225479	47.15100
8	32.97645	21.995024	43.95787
9	29.76466	18.747835	40.78148

Interprete the results. *Your interpretation of the results goes here!*

## Steps of Assignment 3 in detail

Check again the accuracy of the linear regression.

```
In [103]: 1 lm.fit1=lm(medv ~ lstat+rm+nox+dis+ptratio,data=Boston)
          2 summary(lm.fit1)
```

Call:

```
lm(formula = medv ~ lstat + rm + nox + dis + ptratio, data = Boston)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-12.7765	-3.0186	-0.6481	1.9752	27.7625

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	37.49920	4.61295	8.129	3.43e-15 ***
lstat	-0.58108	0.04794	-12.122	< 2e-16 ***
rm	4.16331	0.41203	10.104	< 2e-16 ***
nox	-17.99657	3.26095	-5.519	5.49e-08 ***
dis	-1.18466	0.16842	-7.034	6.64e-12 ***
ptratio	-1.04577	0.11352	-9.212	< 2e-16 ***

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.994 on 500 degrees of freedom

Multiple R-squared: 0.7081, Adjusted R-squared: 0.7052

F-statistic: 242.6 on 5 and 500 DF, p-value: < 2.2e-16

## Add interaction terms

Fit a model with interaction terms. Don't forget to also include the include the plain predictors.

The R syntax `lstat*rm` is a shorthand for `lstat+rm+lstat:rm`, which includes interaction term and plain predictors.

```
In [104]: 1 lm.fit2=lm(medv~lstat*rm+nox+dis+ptratio ,data=Boston)
          2 summary(lm.fit2)
```

Call:

```
lm(formula = medv ~ lstat * rm + nox + dis + ptratio, data = Boston)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-19.3061	-2.4720	-0.3607	1.8192	29.9086

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	3.15175	4.87960	0.646	0.519
lstat	1.81152	0.19612	9.237	< 2e-16 ***
rm	8.33437	0.49109	16.971	< 2e-16 ***
nox	-12.36513	2.88470	-4.286	2.18e-05 ***
dis	-1.01845	0.14776	-6.893	1.66e-11 ***
ptratio	-0.71520	0.10266	-6.967	1.03e-11 ***
lstat:rm	-0.41854	0.03352	-12.488	< 2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.363 on 499 degrees of freedom

Multiple R-squared: 0.7776, Adjusted R-squared: 0.7749

F-statistic: 290.8 on 6 and 499 DF, p-value: < 2.2e-16

Interprete the results. *Your interpretation of the results goes here!*

## Apply non-linear transformations to some predictors

Fit a model with non-linear transformations of the predictor terms. Don't forget to also include the include the plain predictors.

The R the syntax  $I(X^2)$  includes a predictor  $X^2$  but not the plain predictor.

In [113]:

```
1 lm.fit3=lm(medv~lstat*rm+I((lstat*rm)^2)+nox+dis+ptratio,data=Bost
2 summary(lm.fit3)
```

Call:  
lm(formula = medv ~ lstat \* rm + I((lstat \* rm)^2) + nox + dis +  
ptratio, data = Boston)

Residuals:

	Min	1Q	Median	3Q	Max
	-18.4149	-2.4339	-0.2956	1.9426	27.3107

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	1.055e+01	5.499e+00	1.919	0.05558	.
lstat	1.547e+00	2.158e-01	7.167	2.79e-12	***
rm	7.600e+00	5.517e-01	13.777	< 2e-16	***
I((lstat * rm)^2)	3.799e-04	1.335e-04	2.845	0.00462	**
nox	-1.229e+01	2.865e+00	-4.290	2.14e-05	***
dis	-1.064e+00	1.476e-01	-7.209	2.10e-12	***
ptratio	-7.112e-01	1.019e-01	-6.977	9.68e-12	***
lstat:rm	-4.468e-01	3.473e-02	-12.864	< 2e-16	***

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.333 on 498 degrees of freedom  
Multiple R-squared: 0.7812, Adjusted R-squared: 0.7781  
F-statistic: 253.9 on 7 and 498 DF, p-value: < 2.2e-16

The increase of  $R^2$  and the low  $p$ -value associated with the quadratic term suggests that it leads to an improved model. Use ANOVA to check if the quadratic fit is superior to the linear fit.

In R use the `anova()` function to further quantify the extent to which the quadratic fit is superior to the linear fit.

In [114]:

```
1 anova(lm.fit2, lm.fit3)
```

A anova: 2 × 6

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	499	9500.382	NA	NA	NA	NA
2	498	9348.436	1	151.9459	8.094303	0.004623345

Interprete the results. *Your interpretation of the results goes here!*

Check if including additional polynomial terms, up to  $N$  order, lead to an improvement in the model fit.

In R the function `poly(lstat,N)` includes the predictors  $X^1, X^2, \dots, X^N$ .

In [129]:

1

2

lm.fit4=lm(medv~lstat\*rm+poly(lstat,5)+nox+dis+ptratio,data=Boston)

summary(lm.fit4)

Call:  
lm(formula = medv ~ lstat \* rm + poly(lstat, 5) + nox + dis + ptratio, data = Boston)

Residuals:

Min	1Q	Median	3Q	Max
-16.3062	-2.2562	-0.3016	1.8543	28.2698

Coefficients: (1 not defined because of singularities)

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	15.89477	5.93256	2.679	0.007625	**
lstat	1.12040	0.31118	3.601	0.000350	***
rm	6.52913	0.71099	9.183	< 2e-16	***
poly(lstat, 5)1	NA	NA	NA	NA	
poly(lstat, 5)2	17.16543	6.98854	2.456	0.014383	*
poly(lstat, 5)3	-13.81008	4.59170	-3.008	0.002767	**
poly(lstat, 5)4	13.84013	4.53013	3.055	0.002371	**
poly(lstat, 5)5	-15.09359	4.29516	-3.514	0.000482	***
nox	-13.75133	2.82328	-4.871	1.50e-06	***
dis	-1.03260	0.14489	-7.127	3.65e-12	***
ptratio	-0.74069	0.10114	-7.324	9.85e-13	***
lstat:rm	-0.30551	0.05197	-5.878	7.62e-09	***

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.241 on 495 degrees of freedom  
Multiple R-squared: 0.7916, Adjusted R-squared: 0.7873  
F-statistic: 188 on 10 and 495 DF, p-value: < 2.2e-16

Use ANOVA to check if the quadratic fit is superior to the linear fit.

In [123]:

1

anova(lm.fit2, lm.fit4)

A anova: 2 × 6

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	499	9500.382	NA	NA	NA	NA
2	495	8903.772	4	596.6094	8.292038	1.767133e-06

Interpret the results. *Your interpretation of the results goes here!*

Fit and assess other non-linear transofrmations, e.g., log(X).

In [132]:

1

2

3

lm.fit5=lm(medv~poly(lstat,5)+rm+log(rm)+nox+dis+ptratio,data=Bost

summary(lm.fit5)

anova(lm.fit2, lm.fit5)

log(rm)                -137.40385      16.76128      -8.198   2.12e-15   \*\*\*

nox                    -16.64076        2.73369      -6.087   2.30e-09   \*\*\*

dis                    -0.97087        0.14101      -6.885   1.76e-11   \*\*\*

ptratio                -0.78430        0.09663      -8.116   3.83e-15   \*\*\*

----

Signif. codes:   0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.116 on 495 degrees of freedom

Multiple R-squared:   0.8037,      Adjusted R-squared:   0.7997

F-statistic: 202.6 on 10 and 495 DF,   p-value: < 2.2e-16

A anova: 2 × 6

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	499	9500.382	NA	NA	NA	NA
2	495	8386.756	4	1113.626	16.432	1.190966e-12

Interprete the results. *Your interpretation of the results goes here!*

*Beat the teacher!*

# Use categorical predictors

Therefore, we will now examine the Carseats data, which is part of the ISLR library. We will attempt to predict Sales (child car seat sales) in 400 locations based on a number of predictors.

Load the dataset and get an overview of the predictors.

In [135]:

1

load(file = "../ISLR/data/Carseats.rda")

2

summary(Carseats)

Sales		CompPrice		Income		Advertising		
Min.	: 0.000	Min.	: 77	Min.	: 21.00	Min.	: 0.000	
1st Qu.:	5.390	1st Qu.:	115	1st Qu.:	42.75	1st Qu.:	0.000	
Median	: 7.490	Median	:125	Median	: 69.00	Median	: 5.000	
Mean	: 7.496	Mean	:125	Mean	: 68.66	Mean	: 6.635	
3rd Qu.:	9.320	3rd Qu.:	135	3rd Qu.:	91.00	3rd Qu.:	12.000	
Max.	:16.270	Max.	:175	Max.	:120.00	Max.	:29.000	
Population		Price		ShelveLoc		Age		Educ
Min.	: 10.0	Min.	: 24.0	Bad	: 96	Min.	:25.00	Min.
1st Qu.:	139.0	1st Qu.:	100.0	Good	: 85	1st Qu.:	39.75	1st Q
Median	:272.0	Median	:117.0	Medium:	219	Median	:54.50	Median
Mean	:264.8	Mean	:115.8			Mean	:53.32	Mean
3rd Qu.:	398.5	3rd Qu.:	131.0			3rd Qu.:	66.00	3rd Q
Max.	:509.0	Max.	:191.0			Max.	:80.00	Max.
Urban		US						
No	:118	No	:142					
Yes:	282	Yes:	258					

There are categorical predictors `ShelveLoc` , `Urban` , and `US` . Use them in a prediction model.

For the categorical predictors, `R` generates dummy variables.



In [145]:

1

2

lm.fit6=lm(Sales~.,data=Carseats)  
summary(lm.fit6)

Call:  
lm(formula = Sales ~ ., data = Carseats)

Residuals:

Min	1Q	Median	3Q	Max
-2.8692	-0.6908	0.0211	0.6636	3.4115

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	5.6606231	0.6034487	9.380	< 2e-16	***
CompPrice	0.0928153	0.0041477	22.378	< 2e-16	***
Income	0.0158028	0.0018451	8.565	2.58e-16	***
Advertising	0.1230951	0.0111237	11.066	< 2e-16	***
Population	0.0002079	0.0003705	0.561	0.575	
Price	-0.0953579	0.0026711	-35.700	< 2e-16	***
ShelveLocGood	4.8501827	0.1531100	31.678	< 2e-16	***
ShelveLocMedium	1.9567148	0.1261056	15.516	< 2e-16	***
Age	-0.0460452	0.0031817	-14.472	< 2e-16	***
Education	-0.0211018	0.0197205	-1.070	0.285	
UrbanYes	0.1228864	0.1129761	1.088	0.277	
USYes	-0.1840928	0.1498423	-1.229	0.220	

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.019 on 388 degrees of freedom  
Multiple R-squared: 0.8734, Adjusted R-squared: 0.8698  
F-statistic: 243.4 on 11 and 388 DF, p-value: < 2.2e-16

Interprete the results. *Your interpretation of the results goes here!*

Play around with the predictors. For instance:

In [148]:

1

2

lm.fit7=lm(Sales~.-Population-Education-Age-Urban-US +Income:Adver  
summary(lm.fit7)

Call:  
lm(formula = Sales ~ . - Population - Education - Age - Urban -  
US + Income:Advertising + Price:Age, data = Carseats)

Residuals:

Min	1Q	Median	3Q	Max
-2.9433	-0.7721	-0.0059	0.6687	3.7429

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	3.299e+00	4.759e-01	6.932	1.73e-11	***
CompPrice	9.345e-02	4.155e-03	22.492	< 2e-16	***
Income	9.809e-03	2.611e-03	3.756	0.000199	***
Advertising	5.339e-02	2.099e-02	2.544	0.011338	*
Price	-7.590e-02	2.966e-03	-25.591	< 2e-16	***
ShelveLocGood	4.895e+00	1.538e-01	31.825	< 2e-16	***
ShelveLocMedium	1.991e+00	1.265e-01	15.736	< 2e-16	***
Income:Advertising	8.753e-04	2.802e-04	3.124	0.001917	**
Price:Age	-3.682e-04	2.685e-05	-13.713	< 2e-16	***

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.027 on 391 degrees of freedom  
Multiple R-squared: 0.8704, Adjusted R-squared: 0.8677  
F-statistic: 328.2 on 8 and 391 DF, p-value: < 2.2e-16

Beat the teacher!