

Regression

The `library()` function is used to load libraries (functions and data sets that are not included in the base R). Load the `MASS` and `ISLR2` packages.

```
library(MASS)
library(ISLR2)

## 
##     'ISLR2'

## The following object is masked from 'package:MASS':
## 
##     Boston
```

Simple Linear Regression

The `Boston` data: `medv` (median house value) for 506 census tracts in Boston. Predict `medv` using 12 predictors such as `rmvar` (average number of rooms per house), `age` (average age of houses), and `lstat` (percent of households with low socioeconomic status). We will start by using the `lm()` function to fit a simple linear regression model, with `medv` as the response and `lstat` as the predictor.

```
lm.fit <- lm(medv ~ lstat)

## Error in eval(predvars, data, env):     'medv'
```

The command causes an error because R does not know where to find the variables `medv` and `lstat`. The next line tells R that the variables are in `Boston`. If we attach `Boston`, the first line works fine because R now recognizes the variables.

```
lm.fit <- lm(medv ~ lstat, data = Boston)
attach(Boston)
lm.fit <- lm(medv ~ lstat)
```

If we type `lm.fit`, some basic information about the model is output. For more detailed information, we use `summary(lm.fit)`. This gives us *p*-values and standard errors for the coefficients, as well as the R^2 statistic and *F*-statistic for the model.

```
lm.fit

## 
## Call:
## lm(formula = medv ~ lstat)
## 
## Coefficients:
## (Intercept)      lstat
##       34.55      -0.95
summary(lm.fit)

## 
## Call:
## lm(formula = medv ~ lstat)
```

```

## 
## 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

```

To see the list of objects in `lm.fit` and get access to one of it:

```

names(lm.fit)

## [1] "coefficients"   "residuals"      "effects"        "rank"
## [5] "fitted.values"  "assign"         "qr"            "df.residual"
## [9] "xlevels"         "call"          "terms"         "model"

lm.fit$coefficients

## (Intercept)      lstat
## 34.5538409 -0.9500494

```

In order to obtain a confidence interval for the coefficient estimates, we can use the `confint()` command.

```

confint(lm.fit)

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

```

The `predict()` function can be used to produce confidence intervals and prediction intervals for the prediction of `medv` for a given value of `lstat`.

```

predict(lm.fit, data.frame(lstat = c(5, 10, 15)),
       interval = "confidence")

##      fit     lwr     upr
## 1 29.80359 29.00741 30.59978
## 2 25.05335 24.47413 25.63256
## 3 20.30310 19.73159 20.87461

predict(lm.fit, data.frame(lstat = c(5, 10, 15)),
       interval = "prediction")

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

```

Multiple Linear Regression

In order to fit a multiple linear regression model using least squares, we again use the `lm()` function.

```

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

## 
## Call:
## lm(formula = medv ~ lstat + age, data = Boston)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -15.981  -3.978  -1.283   1.968  23.158 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 33.22276   0.73085 45.458 < 2e-16 ***
## lstat        -1.03207   0.04819 -21.416 < 2e-16 ***
## age          0.03454   0.01223   2.826  0.00491 **  
## ---        
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 6.173 on 503 degrees of freedom
## Multiple R-squared:  0.5513, Adjusted R-squared:  0.5495 
## F-statistic: 309 on 2 and 503 DF, p-value: < 2.2e-16

```

Perform a regression using all of the predictors.

```

lm.fit <- lm(medv ~ ., data = Boston)
summary(lm.fit)

## 
## Call:
## lm(formula = medv ~ ., data = Boston)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -15.1304  -2.7673  -0.5814   1.9414  26.2526 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 41.617270  4.936039  8.431 3.79e-16 ***
## crim        -0.121389  0.033000 -3.678 0.000261 *** 
## zn          0.046963  0.013879  3.384 0.000772 *** 
## indus       0.013468  0.062145  0.217 0.828520  
## chas        2.839993  0.870007  3.264 0.001173 **  
## nox         -18.758022 3.851355 -4.870 1.50e-06 *** 
## rm          3.658119  0.420246  8.705 < 2e-16 *** 
## age         0.003611  0.013329  0.271 0.786595  
## dis         -1.490754  0.201623 -7.394 6.17e-13 *** 
## rad         0.289405  0.066908  4.325 1.84e-05 *** 
## tax         -0.012682  0.003801 -3.337 0.000912 *** 
## ptratio     -0.937533  0.132206 -7.091 4.63e-12 *** 
## lstat      -0.552019  0.050659 -10.897 < 2e-16 *** 
## ---        
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 4.798 on 493 degrees of freedom

```

```

## Multiple R-squared:  0.7343, Adjusted R-squared:  0.7278
## F-statistic: 113.5 on 12 and 493 DF,  p-value: < 2.2e-16

We may wish to run a regression excluding age.

lm.fit1 <- lm(medv ~ . - age, data = Boston)
summary(lm.fit1)

## 
## Call:
## lm(formula = medv ~ . - age, data = Boston)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -15.1851  -2.7330  -0.6116   1.8555  26.3838 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 41.525128  4.919684  8.441 3.52e-16 ***
## crim        -0.121426  0.032969 -3.683 0.000256 *** 
## zn          0.046512  0.013766  3.379 0.000785 *** 
## indus       0.013451  0.062086  0.217 0.828577  
## chas        2.852773  0.867912  3.287 0.001085 **  
## nox        -18.485070 3.713714 -4.978 8.91e-07 *** 
## rm          3.681070  0.411230  8.951 < 2e-16 *** 
## dis        -1.506777  0.192570 -7.825 3.12e-14 *** 
## rad         0.287940  0.066627  4.322 1.87e-05 *** 
## tax        -0.012653  0.003796 -3.333 0.000923 *** 
## ptratio     -0.934649  0.131653 -7.099 4.39e-12 *** 
## lstat       -0.547409  0.047669 -11.483 < 2e-16 *** 
## ---        
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 4.794 on 494 degrees of freedom
## Multiple R-squared:  0.7343, Adjusted R-squared:  0.7284
## F-statistic: 124.1 on 11 and 494 DF,  p-value: < 2.2e-16

```

Interaction Terms

The syntax `lstat:black` tells R to include an interaction term between `lstat` and `black`. The syntax `lstat * age` simultaneously includes `lstat`, `age`, and the interaction term `lstat×age` as predictors.

```

summary(lm(medv ~ lstat * age, data = Boston))

## 
## Call:
## lm(formula = medv ~ lstat * age, data = Boston)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -15.806  -4.045  -1.333   2.085  27.552 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 36.0885359  1.4698355 24.553 < 2e-16 ***
## lstat       -1.3921168  0.1674555 -8.313 8.78e-16 ***
```

```

## age          -0.0007209  0.0198792  -0.036   0.9711
## lstat:age    0.0041560  0.0018518   2.244   0.0252 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.149 on 502 degrees of freedom
## Multiple R-squared:  0.5557, Adjusted R-squared:  0.5531
## F-statistic: 209.3 on 3 and 502 DF,  p-value: < 2.2e-16

```

Non-linear Transformations of the Predictors

Given a predictor X , we can create a predictor X^2 using `I(X^2)`. The function `I()` is needed since the `^` has a special meaning in a formula object. We now perform a regression of `medv` onto `lstat` and `lstat^2`.

```

lm.fit2 <- lm(medv ~ lstat + I(lstat^2))
summary(lm.fit2)

##
## Call:
## lm(formula = medv ~ lstat + I(lstat^2))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.2834 -3.8313 -0.5295  2.3095 25.4148
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 42.862007  0.872084  49.15  <2e-16 ***
## lstat       -2.332821  0.123803 -18.84  <2e-16 ***
## I(lstat^2)   0.043547  0.003745  11.63  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.524 on 503 degrees of freedom
## Multiple R-squared:  0.6407, Adjusted R-squared:  0.6393
## F-statistic: 448.5 on 2 and 503 DF,  p-value: < 2.2e-16

```

Compare:

```

summary(lm(medv ~ lstat + lstat^2))

##
## Call:
## lm(formula = medv ~ lstat + lstat^2)
##
## 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 order to create a cubic fit, we can include a predictor of the form $I(X^3)$. A better approach involves using the `poly()` function to create the polynomial within `lm()`.

```

lm.fit5 <- lm(medv ~ poly(lstat, 5))
summary(lm.fit5)

```

```

##
## Call:
## lm(formula = medv ~ poly(lstat, 5))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13.5433 -3.1039 -0.7052  2.0844 27.1153
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 22.5328    0.2318  97.197 < 2e-16 ***
## poly(lstat, 5)1 -152.4595   5.2148 -29.236 < 2e-16 ***
## poly(lstat, 5)2   64.2272   5.2148 12.316 < 2e-16 ***
## poly(lstat, 5)3  -27.0511   5.2148 -5.187 3.10e-07 ***
## poly(lstat, 5)4   25.4517   5.2148  4.881 1.42e-06 ***
## poly(lstat, 5)5  -19.2524   5.2148 -3.692 0.000247 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.215 on 500 degrees of freedom
## Multiple R-squared:  0.6817, Adjusted R-squared:  0.6785
## F-statistic: 214.2 on 5 and 500 DF,  p-value: < 2.2e-16

```

Log transformation:

```
summary(lm(medv ~ log(rm), data = Boston))
```

```

##
## Call:
## lm(formula = medv ~ log(rm), data = Boston)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -19.487 -2.875 -0.104   2.837  39.816
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -76.488     5.028  -15.21 <2e-16 ***
## log(rm)      54.055     2.739   19.73 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.915 on 504 degrees of freedom
## Multiple R-squared:  0.4358, Adjusted R-squared:  0.4347
## F-statistic: 389.3 on 1 and 504 DF,  p-value: < 2.2e-16

```

Qualitative Predictors

We will now examine the `Carseats` data, which is part of the `ISLR2` library. We will attempt to predict `Sales` (child car seat sales) in 400 locations based on a number of predictors. The `Carseats` data includes qualitative predictors such as `Shelveloc`, an indicator of the quality of the shelving location. The predictor `Shelveloc` takes on three possible values: Bad, Medium, and Good. Given a qualitative variable such as `Shelveloc`, R generates dummy variables automatically. Below we fit a multiple regression model that includes some interaction terms.

```
lm.fit <- lm(Sales ~ . + Income:Advertising + Price:Age,
  data = Carseats)
summary(lm.fit)

##
## Call:
## lm(formula = Sales ~ . + Income:Advertising + Price:Age, data = Carseats)
##
## Residuals:
##    Min      1Q  Median      3Q     Max 
## -2.9208 -0.7503  0.0177  0.6754  3.3413 
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)    
## (Intercept)               6.5755654  1.0087470   6.519 2.22e-10 ***
## CompPrice                  0.0929371  0.0041183  22.567 < 2e-16 ***
## Income                      0.0108940  0.0026044   4.183 3.57e-05 ***
## Advertising                 0.0702462  0.0226091   3.107 0.002030 ** 
## Population                  0.0001592  0.0003679   0.433 0.665330  
## Price                     -0.1008064  0.0074399  -13.549 < 2e-16 ***
## ShelveLocGood                4.8486762  0.1528378   31.724 < 2e-16 ***
## ShelveLocMedium              1.9532620  0.1257682   15.531 < 2e-16 ***
## Age                       -0.0579466  0.0159506  -3.633 0.000318 *** 
## Education                   -0.0208525  0.0196131  -1.063 0.288361  
## UrbanYes                    0.1401597  0.1124019   1.247 0.213171  
## USYes                      -0.1575571  0.1489234  -1.058 0.290729  
## Income:Advertising          0.0007510  0.0002784   2.698 0.007290 ** 
## Price:Age                   0.0001068  0.0001333   0.801 0.423812  
## ---                        
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.011 on 386 degrees of freedom
## Multiple R-squared:  0.8761, Adjusted R-squared:  0.8719 
## F-statistic:  210 on 13 and 386 DF,  p-value: < 2.2e-16
```

The `contrasts()` function returns the coding that R uses for the dummy variables. R has created a `ShelveLocGood` dummy variable that takes on a value of 1 if the shelving location is good, and 0 otherwise. It has also created a `ShelveLocMedium` dummy variable that equals 1 if the shelving location is medium, and 0 otherwise. A bad shelving location corresponds to a zero for each of the two dummy variables.

```
attach(Carseats)
contrasts(Shelveloc)

##
##           Good Medium
## Bad          0     0
## Good         1     0
## Medium       0     1
```