

Multiple Linear Regression

Table of Content

- Linear model
 - ▶ Additive
 - ▶ Interpretation
- Model specification
- Model assumptions
- LS estimates
 - ▶ Mean estimates
 - ▶ CI's for each
 - ▶ Predictions for y
- Goodness of fit
- Model diagnosis
- Categorical X's
- Variable selection

1 Case Study: Fuel Efficiency in Automobiles

2 Multiple regression

- Model Specification
- General linear models
- OLS and its Properties
- Compare three models
- Inference
- Confidence Interval
- Prediction Interval
- Model Diagnoses

3 Appendix

Case Study: Fuel Efficiency in Automobiles

Goal of the study: How to build a fuel efficient car?

- Effects of features on a car
- Given a set of features, we'd like to estimate the mean fuel efficiency as well as the efficiency of one car
- Are Asian cars more efficient than cars built in other regions?

Let us answer these questions with Multiple Regression

Data

Dataset: `car_04_regular.csv` , $n = 226$ cars

Feature	Description
Continent	Continent the Car Company is From
Horsepower	Horsepower of the Vehicle
Weight	Weight of the vehicle (thousand lb)
Length	Length of the Vehicle (inches)
Width	Width of the Vehicle (inches)
Seating	Number of Seats in the Vehicle
Cylinders	Number of Engine Cylinders
Displacement	Volume displaced by the cylinders, defined as $\pi/4 \times \text{bore}^2 \times \text{stroke} \times \text{Number of Cylinders}$
Transmission	Type of Transmission (manual, automatic, continuous)

Goal of Study: Rephrased

- 1 Fuel efficiency is measured by Mileage per Gallon, $Y = \text{MPG_City}$
- 2 Predictors: Effects of each feature on Y
- 3 Estimate - the mean MPG_City for all such cars specified below and - predict Y for the particular car described below
- 4 Are cars built by Asian more efficient?
- 5 Investigate the MPG_City for this newly designed American car

Feature	Value
Continent	America
Horsepower	225
Weight	4
Length	180
Width	80
Seating	5
Cylinders	4
Displacement	3.5
Transmission	automatic

A Quick Glimpse at the data

```
data1 <- read.csv("car_04_regular.csv", header=TRUE)
```

```
names(data1)
```

```
## [1] "Make.Model" "Continent" "MPG_City" "MPG_Hwy" "Horsepower"  
## [6] "Weight" "Length" "Width" "Seating" "Cylinders"  
## [11] "Displacement" "Make" "Transmission"
```

```
dim(data1) # 226 cars and 13 variables
```

```
## [1] 226 13
```

A Quick Glimpse at the data

```
str(data1)
```

```
## 'data.frame':    226 obs. of  13 variables:
## $ Make.Model   : chr  "Acura_RL" "Acura_TL" "Acura_TSX" "Acura_RSX" ...
## $ Continent    : chr  "As" "As" "As" "As" ...
## $ MPG_City      : int   18 20 23 25 17 17 20 18 17 16 ...
## $ MPG_Hwy       : int   24 28 32 34 24 23 28 25 24 22 ...
## $ Horsepower    : int  225 270 200 160 252 265 170 220 330 250 ...
## $ Weight        : num   3.9 3.58 3.32 2.77 3.2 ...
## $ Length        : num  197 189 183 172 174 ...
## $ Width         : num  71.6 72.2 69.4 67.9 71.3 77 76.3 76.1 74.6 76.1 ...
## $ Seating       : int   5 5 5 4 2 7 5 5 5 5 ...
## $ Cylinders      : int   6 6 4 4 6 6 4 6 6 6 ...
## $ Displacement: num   3.5 3.2 2.4 2 3 3.5 1.8 3 4.2 2.7 ...
## $ Make          : chr  "Acura" "Acura" "Acura" "Acura" ...
## $ Transmission: chr  "automatic" "automatic" "automatic" "automatic" ...
```


A Quick Glimpse at the data

```
head(data1)
```

```
##   Make.Model Continent MPG_City MPG_Hwy Horsepower Weight Length Width Seating
## 1   Acura_RL        As      18      24      225   3.90   197   71.6      5
## 2   Acura_TL        As      20      28      270   3.58   189   72.2      5
## 3   Acura_TSX        As      23      32      200   3.32   183   69.4      5
## 4   Acura_RSX        As      25      34      160   2.77   172   67.9      4
## 5   Acura_NSX        As      17      24      252   3.20   174   71.3      2
## 6   Acura_MDX        As      17      23      265   4.45   189   77.0      7
##   Cylinders Displacement  Make Transmission
## 1         6           3.5 Acura   automatic
## 2         6           3.2 Acura   automatic
## 3         4           2.4 Acura   automatic
## 4         4           2.0 Acura   automatic
## 5         6           3.0 Acura   automatic
## 6         6           3.5 Acura   automatic
```

1 Case Study: Fuel Efficiency in Automobiles

2 Multiple regression

- Model Specification
- General linear models
- OLS and its Properties
- Compare three models
- Inference
- Confidence Interval
- Prediction Interval
- Model Diagnoses

3 Appendix

Introduction to Multiple regression

Guiding Question: How does Length affect MPG_City?

- It depends on how we model the response. We will investigate **three models** with Length.
- For the ease of presentation, we define some predictors below that we will use in subsequent models:

$x_1 = \text{Length}, \quad x_2 = \text{Horsepower}, \quad x_3 = \text{Width}, \quad x_4 = \text{Seating}$

, $x_5 = \text{Cylinders}, \quad x_6 = \text{Displacement}$

- 1 Case Study: Fuel Efficiency in Automobiles
- 2 Multiple regression
 - Model Specification
 - General linear models
 - OLS and its Properties
 - Compare three models
 - Inference
 - Confidence Interval
 - Prediction Interval
 - Model Diagnoses
- 3 Appendix

Model Specification

M1. Our first model will only contain one predictor, Length:

$$y_i = \beta_0 + \beta_1 \cdot x_{i1} + \epsilon$$

Interpretation of β_1 is that in general, the mean y will change by β_1 if a car is 1" longer. So we can't really peel off the effect of the Length over y .

Additive model: β_1 ?

Model Specification

M2. Next, we add the predictor Horsepower to our model

$$y_i = \beta_0 + \beta_1 \cdot x_{i1} + \beta_2 \cdot x_{i2} + \epsilon$$

Interpretation of β_1 is that in general, the mean y will change by β_1 if a car is 1" longer and the 'Horse Power's' are the same.

Model Specification

M3. Finally, we fit a model with multiple predictors

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \beta_6 x_{i6} + \epsilon$$

Interpretation of β_1 is that in general, the mean y will change by β_1 if a car is 1" longer and the rest of the features are the same.

Model Specification

Question: Are all the β_1 s same in the 3 models above?

No. The effect of Length β_1 depends on the rest of the features in the model!!!!

- 1 Case Study: Fuel Efficiency in Automobiles
- 2 Multiple regression
 - Model Specification
 - General linear models
 - OLS and its Properties
 - Compare three models
 - Inference
 - Confidence Interval
 - Prediction Interval
 - Model Diagnoses
- 3 Appendix

General linear models

In general, We define a multiple regression as

$$Y = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} + \epsilon_i$$

General linear models: Assumptions

Y : response; X_1, X_2, \dots, X_p : explanatory variables

- Linearity Assumption for this model is

$$\mathbf{E}(y_i|x_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

- The homoscedasticity assumption is

$$\mathbf{Var}(y_i|x_{i1}, x_{i2}, \dots, x_{ip}) = \sigma^2$$

- Normality assumption

$$y_i|x_i \stackrel{iid}{\sim} \mathcal{N}(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}, \sigma^2)$$

- 1 Case Study: Fuel Efficiency in Automobiles
- 2 Multiple regression
 - Model Specification
 - General linear models
 - OLS and its Properties
 - Compare three models
 - Inference
 - Confidence Interval
 - Prediction Interval
 - Model Diagnoses
- 3 Appendix

OLS and its Properties

These β parameters are estimated using the same approach as simple regression, specifically by minimizing the sum of squared residuals (RSS):

$$\min_{b_0, b_1, b_2, \dots, b_p} \sum_{i=1}^n (y_i - b_0 - b_1 x_{i1} - b_2 x_{i2} - \dots - b_p x_{ip})^2$$

OLS Estimates

- Each $\hat{\beta}_i$ is normal with mean β_i
- Produce $se(\hat{\beta}_i)$
- z or t interval for β_i based on $\hat{\beta}_i$.

OLS Estimates: Hypothesis Test

To test that

$$\beta_i = 0 \quad \text{vs.} \quad \beta_i = 0$$

which means that given other variables in the model, there is no x_i effect.
We carry out a t-test:

$$t\text{-stat} = \frac{\hat{\beta}_i - 0}{\text{se}(\hat{\beta}_i)}$$

The p-value is:

$$\text{p-value} = 2 \times P(T \text{ variable} > t\text{stat})$$

.

We reject the null hypothesis at an α level if the p-value is $< \alpha$.

OLS Prediction

A 95% Confidence interval for the mean given a set of predictors:

$$\hat{y} \pm 2 \times se(\hat{y})$$

A 95% prediction interval for a future y given a set of predictors:

$$\hat{y} \pm 2 \times \hat{\sigma}.$$

RSS, MSE, RSE

For multiple regression, RSS is estimated as:

$$RSS = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \cdots + \hat{\beta}_p x_{ip}))^2$$

$$MSE = \frac{RSS}{n - p - 1}$$

$$\hat{\sigma} = RSE = \sqrt{MSE}$$

Goodness of Fit: R^2

TSS measures the total variance in the response Y .

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2$$

How much variability is captured in the linear model using this set of predictors? R^2 measures the proportion of variability in Y that can be explained using this set of predictors in the model.

$$R^2 = \frac{TSS - RSS}{TSS}$$

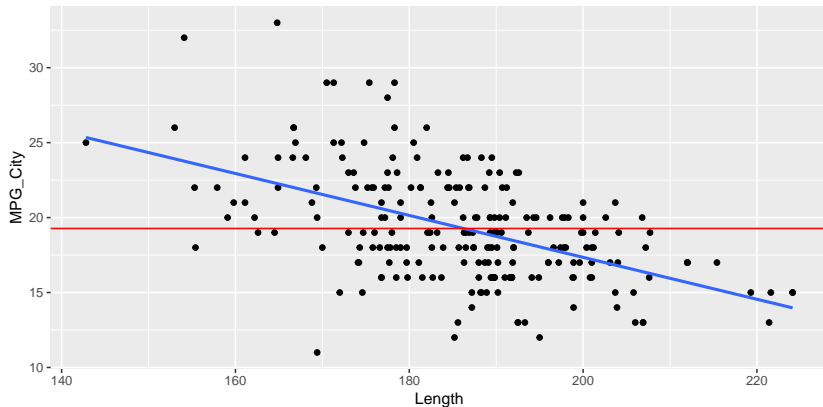
R function 'lm()'

- Linear models are so popular due to their nice interpretations as well as clean solutions
- R-function `lm()` takes a model specification together with other options, outputs all the estimators, summary statistics such as varies sum of squares, standard errors of estimators, testing statistics and p-values.
- The model also outputs the predicted values with margin of errors; confidence intervals and prediction intervals can also be called.

- 1 Case Study: Fuel Efficiency in Automobiles
- 2 Multiple regression
 - Model Specification
 - General linear models
 - OLS and its Properties
 - Compare three models
 - Inference
 - Confidence Interval
 - Prediction Interval
 - Model Diagnoses
- 3 Appendix

Model 1: $\text{MPG_City} \sim \text{Length}$

```
fit1 <- lm(MPG_City ~ Length, data = data1)      # model specification response ~ x1,..
ggplot(data1, aes(x = Length , y = MPG_City)) +
  geom_point() +
  geom_smooth(method="lm", formula = 'y~x', se=F) +
  geom_hline(aes(yintercept = mean(MPG_City)), color = "red")
```



Model 1: MPG_City \sim Length

- We now create a model with `lm()`
- Note from the summary below, the $\hat{\beta}$ for Length is estimated as -0.14.
- We say on average MPG drops .13983 if a car is 1' longer.

```
fit1 <- lm(MPG_City ~ Length, data = data1) # model one
summary(fit1)
```

```
##
## Call:
## lm(formula = MPG_City ~ Length, data = data1)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.626  -2.279  -0.151   1.977  10.731
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  45.3138    2.8975   15.64  <2e-16 ***
## Length       -0.1398    0.0155   -9.01  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.18 on 224 degrees of freedom
## Multiple R-squared:  0.266, Adjusted R-squared:  0.263
## F-statistic: 81.2 on 1 and 224 DF,  p-value: <2e-16
```

Model 2: $\text{MPG_City} \sim \text{Length} + \text{Horsepower}$

```
fit2 <- lm(MPG_City ~ Length + Horsepower, data = data1)
summary(fit2) #sum((fit2$res)^2)
```

```
##
## Call:
## lm(formula = MPG_City ~ Length + Horsepower, data = data1)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.152 -1.558  0.154  1.492  8.563
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  38.62587    2.22525   17.36 < 2e-16 ***
## Length       -0.06191    0.01300    -4.76 3.5e-06 ***
## Horsepower   -0.03690    0.00277   -13.31 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.38 on 223 degrees of freedom
## Multiple R-squared:  0.591, Adjusted R-squared:  0.587
## F-statistic: 161 on 2 and 223 DF, p-value: <2e-16
```

Model 3: Several continuous variables

```
fit3 <- lm(MPG_City ~ Length + Horsepower + Width + Seating +
           Cylinders + Displacement, data = data1)
summary(fit3)

##
## Call:
## lm(formula = MPG_City ~ Length + Horsepower + Width + Seating +
##     Cylinders + Displacement, data = data1)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.877 -1.462  0.073  1.149  8.261
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  45.63372    4.02793   11.33 < 2e-16 ***
## Length       0.04909    0.01730    2.84  0.005 **
## Horsepower   -0.02000    0.00408   -4.90  1.8e-06 ***
## Width        -0.35358    0.06879   -5.14  6.1e-07 ***
## Seating      -0.24135    0.14955   -1.61  0.108
## Cylinders     -0.27169    0.23292   -1.17  0.245
## Displacement -0.93813    0.37166   -2.52  0.012 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.04 on 219 degrees of freedom
## Multiple R-squared:  0.705, Adjusted R-squared:  0.697
## F-statistic: 87.2 on 6 and 219 DF,  p-value: <2e-16
```


Compare 3 Models

Table 3:

	<i>Dependent variable:</i>		
	MPG_City		
	(1)	(2)	(3)
Length	-0.140*** (0.016)	-0.062*** (0.013)	0.049*** (0.017)
Horsepower		-0.037*** (0.003)	-0.020*** (0.004)
Width			-0.354*** (0.069)
Seating			-0.241 (0.150)
Cylinders			-0.272 (0.233)
Displacement			-0.938** (0.372)
Constant	45.300*** (2.900)	38.600*** (2.230)	45.600*** (4.030)
Observations	226	226	226
R ²	0.266	0.591	0.705
Residual Std. Error	3.170 (df = 224)	2.380 (df = 223)	2.040 (df = 219)

Note:

* p<0.1; ** p<0.05; *** p<0.01

Compare 3 Models

- They are different as expected
- Each one has its own meaning!

Question: what does $\hat{\beta}_1$ mean in 3 models

- 1 Case Study: Fuel Efficiency in Automobiles
- 2 Multiple regression
 - Model Specification
 - General linear models
 - OLS and its Properties
 - Compare three models
 - **Inference**
 - Confidence Interval
 - Prediction Interval
 - Model Diagnoses
- 3 Appendix

Questions of interests

- 1 Write down the final OLS equation of `MPG_City` given the rest of the predictors.
- 2 What does each t-interval and t-test do?
- 3 Is `Width` **THE** most important variable, `HP` the second, etc since they each has the smallest p-value,...
- 4 Is `Width` most useful variable due to its largest coefficient in magnitude?
- 5 What is the standard error from the output? Precisely what does it measure?
- 6 Interpret the R^2 reported for this model. Do you feel comfortable using the output for the following questions based on this R^2 value?
- 7 Should we take `Seating` or `Cylinders` out?

- 1 Case Study: Fuel Efficiency in Automobiles
- 2 Multiple regression
 - Model Specification
 - General linear models
 - OLS and its Properties
 - Compare three models
 - Inference
 - **Confidence Interval**
 - Prediction Interval
 - Model Diagnoses
- 3 Appendix

Confidence Interval for the Mean

Base on Model 3, the mean of MPG_City among all cars with the same features as the new design: length=180, HP=225, width=80, seating=5, cylinders=4, displacement=3.5, transmission="automatic", continent="Am" is

$$\begin{aligned}\hat{y} &= 45.63 + 0.05 \times 180 - 0.02 \times 225 - 0.35 \\ &\times 80 - 0.24 \times 5 - 0.27 \times 4 - 0.94 \times 3.5 = 16.17,\end{aligned}$$

Confidence Interval for the Mean

```
predict(fit3, newcar, interval = "confidence", se.fit = TRUE)
```

```
## $fit  
##      fit   lwr   upr  
## 1 16.1 14.6 17.7  
##  
## $se.fit  
## [1] 0.784  
##  
## $df  
## [1] 219  
##  
## $residual.scale  
## [1] 2.04
```

Q: What assumptions are needed to make this a valid confidence interval?

- 1 Case Study: Fuel Efficiency in Automobiles
- 2 Multiple regression
 - Model Specification
 - General linear models
 - OLS and its Properties
 - Compare three models
 - Inference
 - Confidence Interval
 - Prediction Interval
 - Model Diagnoses
- 3 Appendix

Prediction Interval

Base on Model 3, MPG_City for this particular new design is

$$\begin{aligned}\hat{y} &= 45.63 + 0.05 \times 180 - 0.02 \times 225 - 0.35 \\ &\times 80 - 0.24 \times 5 - 0.27 \times 4 - 0.94 \times 3.5 = 16.17\end{aligned}$$

with a 95% prediction interval approximately to be

$$\hat{y} \pm 2 \times RSE = 16.17 \pm 2 \times 2.036.$$

Prediction Interval

```
# future prediction intervals
predict(fit3, newcar, interval = "predict", se.fit = TRUE)

## $fit
##      fit   lwr   upr
## 1 16.1 11.8 20.4
##
## $se.fit
## [1] 0.784
##
## $df
## [1] 219
##
## $residual.scale
## [1] 2.04
```

Q: What assumptions are needed to make this a valid prediction interval?

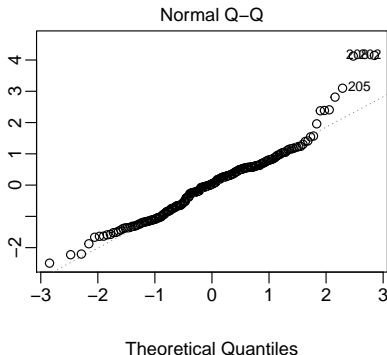
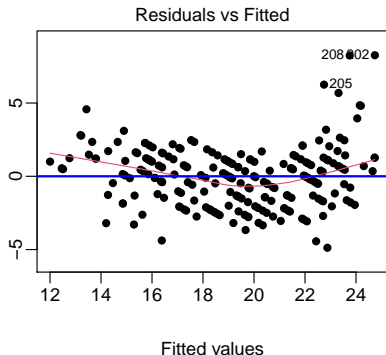
- 1 Case Study: Fuel Efficiency in Automobiles
- 2 Multiple regression
 - Model Specification
 - General linear models
 - OLS and its Properties
 - Compare three models
 - Inference
 - Confidence Interval
 - Prediction Interval
 - Model Diagnoses
- 3 Appendix

Model Diagnoses

To check the model assumptions are met, we examine the residual plot and the qqplot of the residuals.

We use the first and second plots of `plot(fit)`.

```
par(mfrow=c(1,2), mar=c(5,2,4,2), mgp=c(3,0.5,0)) # plot(fit3) produces several plots
plot(fit3, 1, pch=16) # residual plot
abline(h=0, col="blue", lwd=2)
plot(fit3, 2) # qqplot
```



Model Diagnoses

Are the linear model assumptions met for the model fit here (fit3)? What might be violated?

- linearity?
- Equal variances?
- Normality?

- 1 Case Study: Fuel Efficiency in Automobiles
- 2 Multiple regression
 - Model Specification
 - General linear models
 - OLS and its Properties
 - Compare three models
 - Inference
 - Confidence Interval
 - Prediction Interval
 - Model Diagnoses
- 3 Appendix

Goodness of Fit: R^2

Remark 1:

- $TSS \geq RSS$. Why so???
- $R^2 \leq 1$.
- $TSS = RSS + \sum(\hat{y}_i - \bar{y})^2$.
- $(\text{corr}(y, \hat{y}))^2$

An R^2 statistic that is close to 1 indicates that a large proportion of the variability in the response has been explained by the regression.

Goodness of Fit: R^2

Remark2:

- How large R^2 needs to be so that you are comfortable to use the linear model?
- Though R^2 is a very popular notion of goodness of fit, but it has its limitation. Mainly all the sum of squared errors defined so far are termed as Training Errors. It really only measures how good a model fits the data that we use to build the model. It may not generate well to unseen data.