

Complicating the RHS of the Linear Model

Fundamental Techniques in Data Science with R



**Utrecht
University**

Kyle M. Lang

Department of Methodology & Statistics
Utrecht University

Outline

Categorical Predictors

- Dummy Coding

- Significance Testing for Dummy Codes

Moderation

- Categorical Moderators

Polynomial Regression



Categorical Predictors

Most of the predictors we've considered thus far have been *quantitative*.

- Continuous variables that can take any real value in their range
- Interval or Ratio scaling

We often want to include grouping factors as predictors.

- These variables are *qualitative*.
 - Their values are simply labels.
 - There is no ordering of the categories.
 - Nominal scaling



How to Model Categorical Predictors

We need to be careful when we include categorical predictors into a regression model.

- The variables need to be coded before entering the model

Consider the following indicator of major:

$$X_{maj} = \{1 = Law, 2 = Economics, 3 = Data Science\}$$

- What would happen if we naïvely used this variable to predict program satisfaction?



How to Model Categorical Predictors

```
mDat <- readRDS("../data/major_data.rds")
```

```
mDat[seq(25, 150, 25), ]
```

	sat	majF	majN
25	1.9	law	1
50	1.4	law	1
75	4.3	econ	2
100	4.1	econ	2
125	5.7	ds	3
150	5.1	ds	3

```
out1 <- lm(sat ~ majN, data = mDat)
```

How to Model Categorical Predictors

```
partSummary(out1, -1)
```

Residuals:

Min	1Q	Median	3Q	Max
-1.303	-0.313	-0.113	0.342	1.342

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.33200	0.12060	-2.753	0.00664
majN	2.04500	0.05582	36.632	< 2e-16

Residual standard error: 0.5582 on 148 degrees of freedom

Multiple R-squared: 0.9007, Adjusted R-squared: 0.9

F-statistic: 1342 on 1 and 148 DF, p-value: < 2.2e-16

Dummy Coding

The most common way to code categorical predictors is *dummy coding*.

- A G -level factor must be converted into a set of $G - 1$ dummy codes.
- Each code is a variable on the dataset that equals 1 for observations corresponding to the code's group and equals 0, otherwise.
- The group without a code is called the *reference group*.



Example Dummy Code

Let's look at the simple example of coding biological sex:

	sex	male
1	female	0
2	male	1
3	male	1
4	female	0
5	male	1
6	female	0
7	female	0
8	male	1
9	female	0
10	female	0



Example Dummy Codes

Now, a slightly more complex example:

	drink	juice	tea
1	juice	1	0
2	coffee	0	0
3	tea	0	1
4	tea	0	1
5	tea	0	1
6	tea	0	1
7	juice	1	0
8	tea	0	1
9	coffee	0	0
10	juice	1	0



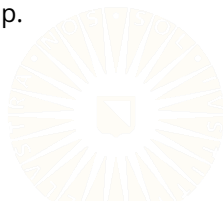
Using Dummy Codes

To use the dummy codes, we simply include the $G - 1$ codes as $G - 1$ predictor variables in our regression model.

$$Y = \beta_0 + \beta_1 X_{male} + \varepsilon$$

$$Y = \beta_0 + \beta_1 X_{juice} + \beta_2 X_{tea} + \varepsilon$$

- The intercept corresponds to the mean of Y for the reference group.
- Each slope represents the difference between the mean of Y in the coded group and the mean of Y in the reference group.



Example

First, an example with a single, binary dummy code:

```
## Read in some data:  
cDat <- readRDS("../data/cars_data.rds")  
  
## Fit and summarize the model:  
out2 <- lm(price ~ mtOpt, data = cDat)
```

Example

```
partSummary(out2, -1)
```

Residuals:

Min	1Q	Median	3Q	Max
-10.341	-6.338	-3.141	2.662	38.059

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	23.841	1.623	14.691	<2e-16
mtOpt	-6.603	2.004	-3.295	0.0014

Residual standard error: 9.18 on 91 degrees of freedom

Multiple R-squared: 0.1066, Adjusted R-squared: 0.09679

F-statistic: 10.86 on 1 and 91 DF, p-value: 0.001403

Interpretations

- The average price of a car without the option for a manual transmission is $\hat{\beta}_0 = 23.84$ thousand dollars.
- The average difference in price between cars that have manual transmissions as an option and those that do not is $\hat{\beta}_1 = -6.6$ thousand dollars.



Example

Fit a more complex model:

```
out3 <- lm(price ~ front + rear, data = cDat)
partSummary(out3, -1)
```

Residuals:

Min	1Q	Median	3Q	Max
-14.050	-6.250	-1.236	3.264	32.950

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	17.63000	2.76119	6.385	7.33e-09
front	-0.09418	2.96008	-0.032	0.97469
rear	11.32000	3.51984	3.216	0.00181

Residual standard error: 8.732 on 90 degrees of freedom

Multiple R-squared: 0.2006, Adjusted R-squared: 0.1829

F-statistic: 11.29 on 2 and 90 DF, p-value: 4.202e-05

Interpretations

- The average price of a four-wheel-drive car is $\hat{\beta}_0 = 17.63$ thousand dollars.
- The average difference in price between front-wheel-drive cars and four-wheel-drive cars is $\hat{\beta}_1 = -0.09$ thousand dollars.
- The average difference in price between rear-wheel-drive cars and four-wheel-drive cars is $\hat{\beta}_2 = 11.32$ thousand dollars.



Example

Include two sets of dummy codes:

```
out4 <- lm(price ~ mtOpt + front + rear, data = cDat)
partSummary(out4, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	21.7187	2.9222	7.432	6.25e-11
mtOpt	-5.8410	1.8223	-3.205	0.00187
front	-0.2598	2.8189	-0.092	0.92677
rear	10.5169	3.3608	3.129	0.00237

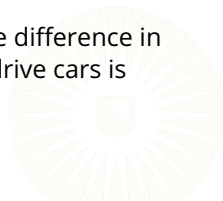
Residual standard error: 8.314 on 89 degrees of freedom

Multiple R-squared: 0.2834, Adjusted R-squared: 0.2592

F-statistic: 11.73 on 3 and 89 DF, p-value: 1.51e-06

Interpretations

- The average price of a four-wheel-drive car that does not have a manual transmission option is $\hat{\beta}_0 = 21.72$ thousand dollars.
- After controlling for drive type, the average difference in price between cars that have manual transmissions as an option and those that do not is $\hat{\beta}_1 = -5.84$ thousand dollars.
- After controlling for transmission options, the average difference in price between front-wheel-drive cars and four-wheel-drive cars is $\hat{\beta}_2 = -0.26$ thousand dollars.
- After controlling for transmission options, the average difference in price between rear-wheel-drive cars and four-wheel-drive cars is $\hat{\beta}_3 = 10.52$ thousand dollars.



Significance Testing

For variables with only two levels, we can test the overall factor's significance by evaluating the significance of a single dummy code.

```
partSummary(out2, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	23.841	1.623	14.691	<2e-16
mtOpt	-6.603	2.004	-3.295	0.0014

Residual standard error: 9.18 on 91 degrees of freedom

Multiple R-squared: 0.1066, Adjusted R-squared: 0.09679

F-statistic: 10.86 on 1 and 91 DF, p-value: 0.001403

Significance Testing

For variables with more than two levels, we need to simultaneously evaluate the significance of each of the variable's dummy codes.

```
partSummary(out4, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	21.7187	2.9222	7.432	6.25e-11
mtOpt	-5.8410	1.8223	-3.205	0.00187
front	-0.2598	2.8189	-0.092	0.92677
rear	10.5169	3.3608	3.129	0.00237

Residual standard error: 8.314 on 89 degrees of freedom

Multiple R-squared: 0.2834, Adjusted R-squared: 0.2592

F-statistic: 11.73 on 3 and 89 DF, p-value: 1.51e-06

Significance Testing

```
summary(out4)$r.squared - summary(out2)$r.squared
```

```
[1] 0.1767569
```

```
anova(out2, out4)
```

Analysis of Variance Table

Model 1: price ~ mtOpt

Model 2: price ~ mtOpt + front + rear

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	91	7668.9				
2	89	6151.6	2	1517.3	10.976	5.488e-05 ***

Signif. codes:

0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Significance Testing

For models with a single nominal factor is the only predictor, we use the omnibus F-test.

```
partSummary(out3, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	17.63000	2.76119	6.385	7.33e-09
front	-0.09418	2.96008	-0.032	0.97469
rear	11.32000	3.51984	3.216	0.00181

Residual standard error: 8.732 on 90 degrees of freedom

Multiple R-squared: 0.2006, Adjusted R-squared: 0.1829

F-statistic: 11.29 on 2 and 90 DF, p-value: 4.202e-05

MODERATION



Moderation

So far we've been discussing *additive models*.

- Additive models allow us to examine the partial effects of several predictors on some outcome.
 - The effect of one predictor does not change based on the values of other predictors.

Now, we'll discuss *moderation*.

- Moderation allows us to ask *when* one variable, X , affects another variable, Y .
 - We're considering the conditional effects of X on Y given certain levels of a third variable Z .



Equations

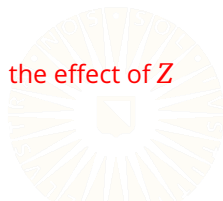
In additive MLR, we might have the following equation:

$$Y = \beta_0 + \beta_1 X + \beta_2 Z + \varepsilon$$

This additive equation assumes that X and Z are independent predictors of Y .

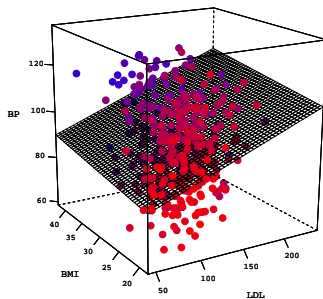
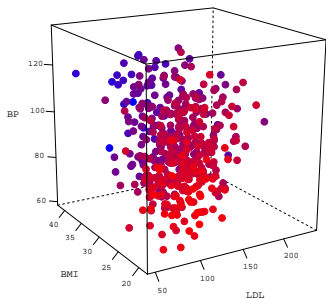
When X and Z are independent predictors, the following are true:

- X and Z *can* be correlated.
- β_1 and β_2 are *partial* regression coefficients.
- The effect of X on Y is the same at **all levels** of Z , and the effect of Z on Y is the same at **all levels** of X .



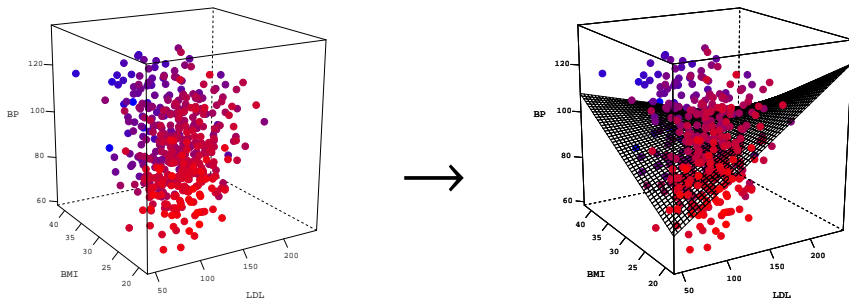
Additive Regression

The effect of X on Y is the same at **all levels** of Z .



Moderated Regression

The effect of X on Y varies **as a function** of Z .



Equations

The following derivation is adapted from Hayes (2017).

- When testing moderation, we hypothesize that the effect of X on Y varies as a function of Z .
- We can represent this concept with the following equation:

$$Y = \beta_0 + f(Z)X + \beta_2Z + \varepsilon \quad (1)$$



Equations

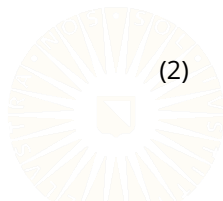
The following derivation is adapted from Hayes (2017).

- When testing moderation, we hypothesize that the effect of X on Y varies as a function of Z .
- We can represent this concept with the following equation:

$$Y = \beta_0 + f(Z)X + \beta_2Z + \varepsilon \quad (1)$$

- If we assume that Z linearly (and deterministically) affects the relationship between X and Y , then we can take:

$$f(Z) = \beta_1 + \beta_3Z \quad (2)$$



Equations

- Substituting Equation 2 into Equation 1 leads to:

$$Y = \beta_0 + (\beta_1 + \beta_3 Z)X + \beta_2 Z + \varepsilon$$



Equations

- Substituting Equation 2 into Equation 1 leads to:

$$Y = \beta_0 + (\beta_1 + \beta_3 Z)X + \beta_2 Z + \varepsilon$$

- Which, after distributing X and reordering terms, becomes:

$$Y = \beta_0 + \beta_1 X + \beta_2 Z + \beta_3 XZ + \varepsilon$$



Testing Moderation

Now, we have an estimable regression model that quantifies the linear moderation we hypothesized.

$$Y = \beta_0 + \beta_1 X + \beta_2 Z + \beta_3 XZ + \varepsilon$$

- To test for significant moderation, we simply need to test the significance of the interaction term, XZ .
 - Check if $\hat{\beta}_3$ is significantly different from zero.

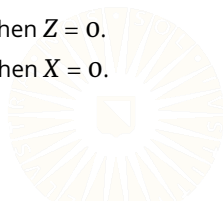


Interpretation

Given the following equation:

$$Y = \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 Z + \hat{\beta}_3 XZ + \hat{\varepsilon}$$

- $\hat{\beta}_3$ quantifies the effect of Z on the focal effect (the $X \rightarrow Y$ effect).
 - For a unit change in Z , $\hat{\beta}_3$ is the expected change in the effect of X on Y .
- $\hat{\beta}_1$ and $\hat{\beta}_2$ are *conditional effects*.
 - Interpreted where the other predictor is zero.
 - For a unit change in X , $\hat{\beta}_1$ is the expected change in Y , when $Z = 0$.
 - For a unit change in Z , $\hat{\beta}_2$ is the expected change in Y , when $X = 0$.



Example

Still looking at the *diabetes* dataset.

- We suspect that patients' BMIs are predictive of their average blood pressure.
- We further suspect that this effect may be differentially expressed depending on the patients' LDL levels.



Example

```
## Focal Effect:
```

```
out0 <- lm(bp ~ bmi, data = dDat)
```

```
partSummary(out0, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	61.9973	3.6659	16.91	<2e-16
bmi	1.2379	0.1371	9.03	<2e-16

Residual standard error: 12.72 on 440 degrees of freedom

Multiple R-squared: 0.1563, Adjusted R-squared: 0.1544

F-statistic: 81.54 on 1 and 440 DF, p-value: < 2.2e-16

Example

```
## Additive Model:
```

```
out1 <- lm(bp ~ bmi + ldl, data = dDat)  
partSummary(out1, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	59.26577	3.91281	15.147	< 2e-16
bmi	1.16567	0.14156	8.235	2.08e-15
ldl	0.04016	0.02056	1.953	0.0515

Residual standard error: 12.68 on 439 degrees of freedom

Multiple R-squared: 0.1636, Adjusted R-squared: 0.1598

F-statistic: 42.94 on 2 and 439 DF, p-value: < 2.2e-16

Example

```
## Moderated Model:
```

```
out2 <- lm(bp ~ bmi * ldl, data = dDat)
partSummary(out2, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	14.480616	14.291677	1.013	0.311514
bmi	2.867825	0.541312	5.298	1.86e-07
ldl	0.448771	0.127160	3.529	0.000461
bmi:ldl	-0.015352	0.004716	-3.255	0.001221

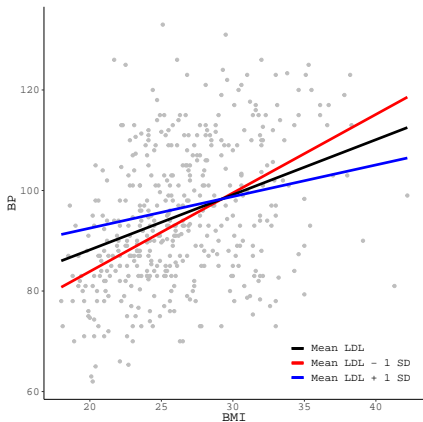
Residual standard error: 12.54 on 438 degrees of freedom

Multiple R-squared: 0.1834, Adjusted R-squared: 0.1778

F-statistic: 32.78 on 3 and 438 DF, p-value: < 2.2e-16

Visualizing the Interaction

We can get a better idea of the patterns of moderation by plotting the focal effect at conditional values of the moderator.



Categorical Moderators

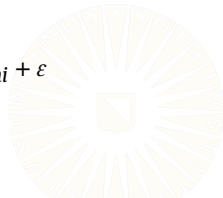
Categorical moderators encode *group-specific* effects.

- E.g., if we include *sex* as a moderator, we are modeling separate focal effects for males and females.

Given a set of codes representing our moderator, we specify the interactions as before:

$$Y_{total} = \beta_0 + \beta_1 X_{inten} + \beta_2 Z_{male} + \beta_3 X_{inten} Z_{male} + \varepsilon$$

$$Y_{total} = \beta_0 + \beta_1 X_{inten} + \beta_2 Z_{lo} + \beta_3 Z_{mid} + \beta_4 Z_{hi} \\ + \beta_5 X_{inten} Z_{lo} + \beta_6 X_{inten} Z_{mid} + \beta_7 X_{inten} Z_{hi} + \varepsilon$$



Example

```
## Load data:
socSup <- readRDS(paste0(dataDir, "social_support.rds"))

## Focal effect:
out3 <- lm(bdi ~ tanSat, data = socSup)
partSummary(out3, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	24.4089	5.3502	4.562	1.54e-05
tanSat	-0.8100	0.3124	-2.593	0.0111

Residual standard error: 9.278 on 93 degrees of freedom

Multiple R-squared: 0.06742, Adjusted R-squared: 0.05739

F-statistic: 6.723 on 1 and 93 DF, p-value: 0.01105

Example

```
## Estimate the interaction:
```

```
out4 <- lm(bdi ~ tanSat * sex, data = socSup)
partSummary(out4, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	20.8478	6.2114	3.356	0.00115
tanSat	-0.5772	0.3614	-1.597	0.11372
sexmale	14.3667	12.2054	1.177	0.24223
tanSat:sexmale	-0.9482	0.7177	-1.321	0.18978

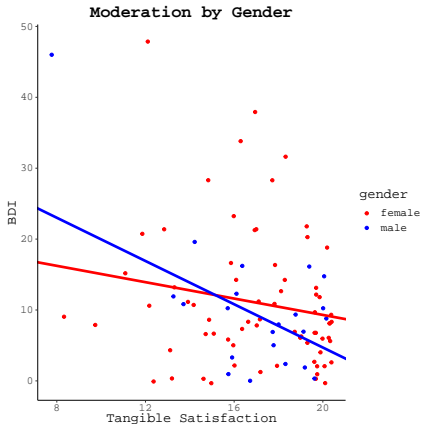
Residual standard error: 9.267 on 91 degrees of freedom

Multiple R-squared: 0.08955, Adjusted R-squared: 0.05954

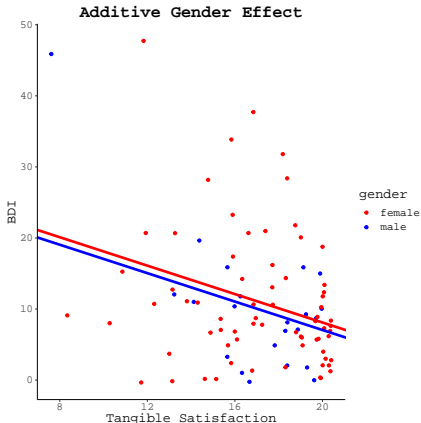
F-statistic: 2.984 on 3 and 91 DF, p-value: 0.03537

Visualizing Categorical Moderation

$$\hat{Y}_{BDI} = 20.85 - 0.58X_{tsat} + 14.37Z_{male} - 0.95X_{tsat}Z_{male}$$



$$\hat{Y}_{BDI} = 28.10 - 1.00X_{tsat} - 1.05Z_{male}$$



POLYNOMIAL REGRESSION



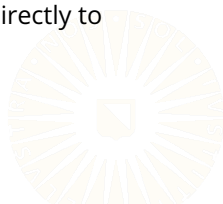
Polynomial Regression

Polynomial regression simply incorporates powered transformations of the predictors into the model.

- Polynomial terms (i.e., power terms) model curvature in the relationships.

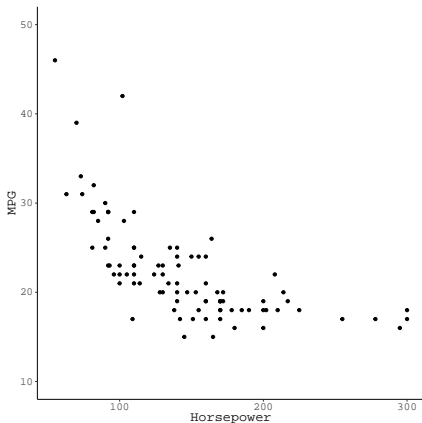
We can think about polynomial terms as interactions between a predictor and itself.

- Many of the rules that apply to interactions transfer directly to polynomials.



Polynomial Visualization

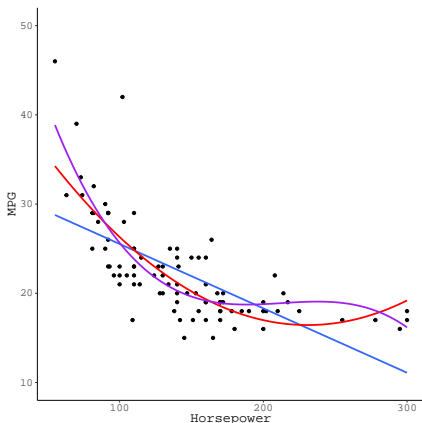
We may hypothesize a curvilinear relationship between X and Y .



Polynomial Visualization

Polynomials are one way to model curvilinear relationships.

- $\hat{Y}_{mpg} = \hat{\beta}_0 + \hat{\beta}_1 X_{hp}$
- $\hat{Y}_{mpg} = \hat{\beta}_0 + \hat{\beta}_1 X_{hp} + \hat{\beta}_2 X_{hp}^2$
- $\hat{Y}_{mpg} = \hat{\beta}_0 + \hat{\beta}_1 X_{hp} + \hat{\beta}_2 X_{hp}^2 + \hat{\beta}_3 X_{hp}^3$



Example

```
## Attach the data:
data(Cars93)

## Fit the linear model:
out6 <- lm(MPG.city ~ Horsepower, data = Cars93)

## Fit the quadratic model:
out7 <- lm(MPG.city ~ Horsepower + I(Horsepower^2),
           data = Cars93)
```

Example

```
partSummary(out6, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	32.746279	1.273229	25.719	< 2e-16
Horsepower	-0.072174	0.008323	-8.671	1.54e-13

Residual standard error: 4.181 on 91 degrees of freedom

Multiple R-squared: 0.4524, Adjusted R-squared: 0.4464

F-statistic: 75.19 on 1 and 91 DF, p-value: 1.537e-13

- For each unit increase in horsepower, the expected change in fuel economy is $\hat{\beta}_1 = -0.0722$ units.

Example

```
partSummary(out7, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.714e+01	2.544e+00	18.528	< 2e-16
Horsepower	-2.660e-01	3.186e-02	-8.350	7.71e-13
I(Horsepower^2)	5.762e-04	9.239e-05	6.237	1.42e-08

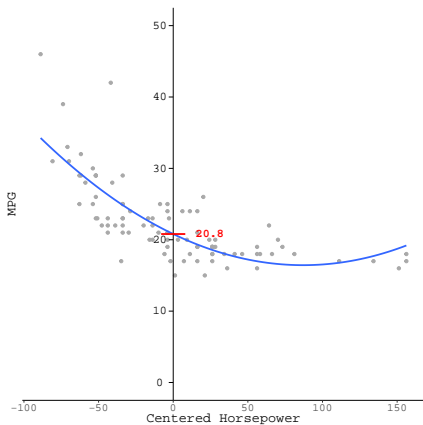
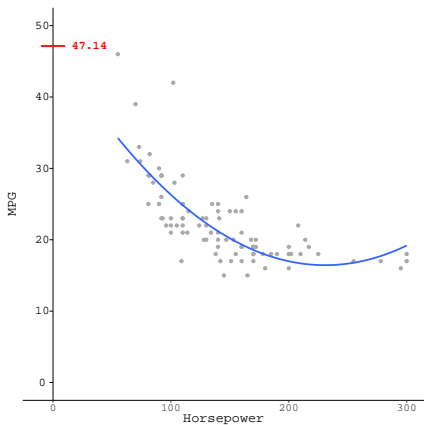
Residual standard error: 3.513 on 90 degrees of freedom

Multiple R-squared: 0.6177, Adjusted R-squared: 0.6092

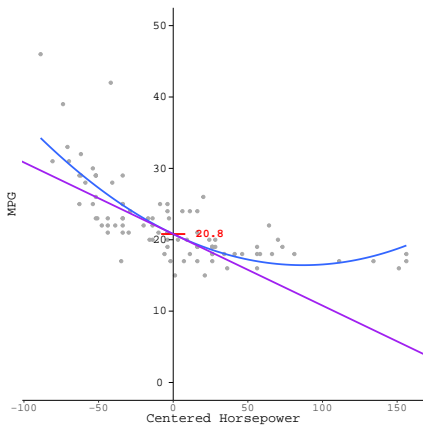
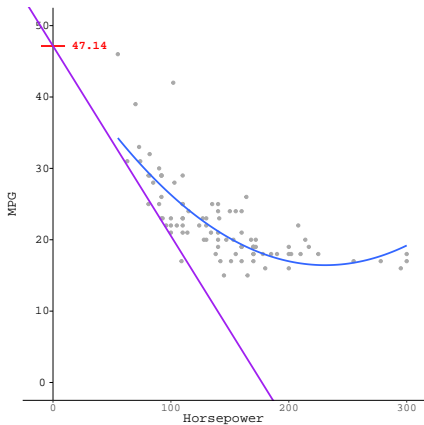
F-statistic: 72.7 on 2 and 90 DF, p-value: < 2.2e-16

- Extrapolating from powerless cars, each unit increase in horsepower, is expected to change fuel economy by $\hat{\beta}_1 = -0.266$ units.
- For a unit increase in horsepower, the effect of horsepower on fuel economy is expected to increase by $\hat{\beta}_2 = 5.76 \times 10^{-4}$ units.

Effects of Centering



Effects of Centering



Example

```
## Mean center horsepower:
```

```
Cars93 <- mutate(Cars93, HorsepowerMC = Horsepower - mean(Horsepower))
```

```
## Fit the quadratic model:
```

```
out8 <- lm(MPG.city ~ HorsepowerMC + I(HorsepowerMC^2), data = Cars93)
```



Example

```
partSummary(out8, -c(1, 2))
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.080e+01	4.422e-01	47.038	< 2e-16
HorsepowerMC	-1.003e-01	8.320e-03	-12.053	< 2e-16
I(HorsepowerMC^2)	5.762e-04	9.239e-05	6.237	1.42e-08

Residual standard error: 3.513 on 90 degrees of freedom

Multiple R-squared: 0.6177, Adjusted R-squared: 0.6092

F-statistic: 72.7 on 2 and 90 DF, p-value: < 2.2e-16

- Averaging over cars, each unit increase in horsepower, is expected to change fuel economy by $\hat{\beta}_1 = -0.1003$ units.
- For a unit increase in horsepower, the effect of horsepower on fuel economy is expected to increase by $\hat{\beta}_2 = 5.76 \times 10^{-4}$ units.

References

Hayes, A. F. (2017). *Introduction to mediation, moderation, and conditional process analysis: A regression-based approach*. New York: Guilford Press.

