

# Quantitative Predictors and Interactions

Prof. Lauren Perry

## Qualitative Predictors

So far, we've focused on only quantitative predictors.

Often, data sets have one or more *qualitative* predictors.

We need to consider how to fit these into a numeric model fitting context.

## Qualitative Predictors with Two Levels

Consider the variable `Own` from the `Credit` data.

```
credit <- read.csv("C:/Users/cappiello/OneDrive - California State Univers  
own <- as.factor(credit$Own)  
summary(own)
```

```
## No Yes
```

```
## 193 207
```

To put this into a regression model, we use a *dummy variable*:

$$x_i = I(\text{the } i\text{th person owns a house})$$

## Qualitative Predictors with Two Levels

This results in the model

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

which takes values

- ▶  $\beta_0 + \beta_1 + \epsilon_i$  if the  $i$ th person owns a house.

and

- ▶  $\beta_0 + \epsilon_i$  if the  $i$ th person does not own a house.

So  $\beta_1$  is the average difference in credit card balance between owners and non-owners.

## Qualitative Predictors with Two Levels

```
summary(lm(Limit ~ Own, data = credit))
```

Coefficients: Estimate Std. Error t value Pr(>|t|)

(Intercept) 4713.17 166.35 28.333 <2e-16 \*\*\*

OwnYes 43.35 231.24 0.187 0.851

---

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

Residual standard error: 2311 on 398 degrees of freedom

Multiple R-squared: 8.83e-05, Adjusted R-squared: -0.002424

F-statistic: 0.03515 on 1 and 398 DF, p-value: 0.8514

## Qualitative Predictors with More than Two Levels

Consider the variable `region` from the `Credit` data.

```
##   East South  West
##    99   199   102
```

We can represent this by constructing *two* dummy variables.

$$x_{i,1} = I(\text{ith person is from the South})$$

$$x_{i,2} = I(\text{ith person is from the West})$$

## Qualitative Predictors with More than Two Levels

Using region to predict credit,

$$y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \epsilon_i$$

Why only two dummy variables? Consider:

- ▶ If the  $i$ th person is from the South,  $y_i = \beta_0 + \beta_1 x_{i,1} + \epsilon_i$ .
- ▶ If the  $i$ th person is from the West,  $y_i = \beta_0 + \beta_2 x_{i,2} + \epsilon_i$
- ▶ If the  $i$ th person is from the East,  $y_i = \beta_0 + \epsilon_i$

So each factor is represented in the model.

- ▶ Because East has no dummy variable, it is known as the *baseline*.

## Qualitative Predictors with More than Two Levels

```
summary(lm(Limit ~ Region, data = credit))
```

Coefficients: Estimate Std. Error t value Pr(>|t|)

(Intercept) 4881.6 232.4 21.009 <2e-16 \*\*\*

RegionSouth -153.1 284.3 -0.539 0.590

RegionWest -273.8 326.2 -0.839 0.402

---

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

Residual standard error: 2312 on 397 degrees of freedom

Multiple R-squared: 0.001781, Adjusted R-squared: -0.003248

F-statistic: 0.3541 on 2 and 397 DF, p-value: 0.702



## Qualitative Predictors

We can also use this approach for a mix of qualitative and quantitative variables in a model.

```
mod2 <- lm(Limit ~ Income + Rating + Own + Region, data=credit)
summary(mod2)
```

```
Coefficients: Estimate Std. Error t value Pr(>|t|)
(Intercept) -539.62205    30.68155  -17.588  <2e-16 ***
Income       0.55281     0.42508   1.300   0.194
Rating      14.77373     0.09685 152.545  <2e-16 ***
OwnYes       2.78064     18.30426  0.152   0.879
RegionSouth  0.71509     22.49522  0.032   0.975
RegionWest  18.21038     25.82151  0.705   0.481
```

# Accounting for Interactions

Sometimes, two predictor variables *interact* in their impact on the outcome.

Example:

- ▶ Suppose spending money on TV advertising increases the effectiveness of radio advertising.
- ▶ We want a way to let  $\beta_{\text{radio}}$  vary based on values of TV...

## Accounting for Interactions

Consider

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon$$

How does this let  $\beta_{\text{radio}}$  vary based on values of  $X_2 = \text{TV}$ ?

$$\begin{aligned} Y &= \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2 + \epsilon \\ &= \beta_0 + \tilde{\beta}_1 X_1 + \beta_2 X_2 + \epsilon \end{aligned}$$

We can interpret  $\beta_3$  as the increase in effectiveness of TV advertising associated with a one-unit increase in radio advertising (or vice versa).

Consider: Why does estimating the coefficients not require any changes to our least squares approach?

```
Advertising <- read.csv("~/Courses/STAT 196M/datasets/Advertising.csv")
mod3 <- lm(sales ~ TV + radio + TV*radio, data=Advertising)
summary(mod3)
```

```
Coefficients: Estimate Std. Error t value Pr(>|t|)
(Intercept)  6.750e+00  2.479e-01  27.233  <2e-16 ***
TV           1.910e-02  1.504e-03  12.699  <2e-16 ***
radio       2.886e-02  8.905e-03   3.241  0.0014 **
TV:radio    1.086e-03  5.242e-05  20.727  <2e-16 ***
```

---

Residual standard error: 0.9435 on 196 degrees of freedom

Multiple R-squared: 0.9678, Adjusted R-squared: 0.9673

F-statistic: 1963 on 3 and 196 DF, p-value: < 2.2e-16

Consider  $R_{adj}^2$  for the main effects model:

```
Advertising <- read.csv("~/Courses/STAT 196M/datasets/Advertising.csv")
mod4 <- lm(sales ~ TV + radio, data=Advertising)
summary(mod4)
```

```
Coefficients: Estimate Std. Error t value Pr(>|t|)
(Intercept)  2.92110    0.29449   9.919  <2e-16 ***
TV           0.04575    0.00139  32.909  <2e-16 ***
radio       0.18799    0.00804  23.382  <2e-16 ***
```

---

Residual standard error: 1.681 on 197 degrees of freedom

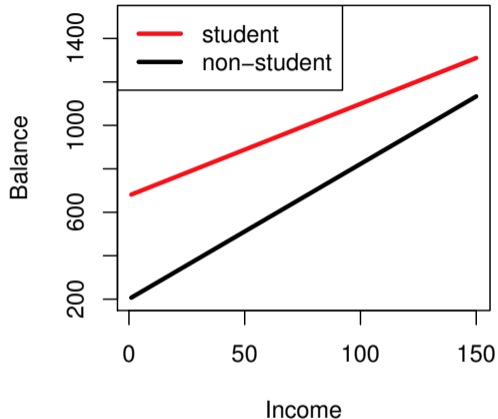
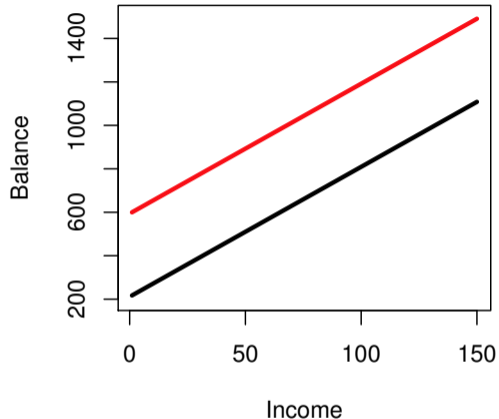
Multiple R-squared: 0.8972, Adjusted R-squared: 0.8962

F-statistic: 859.6 on 2 and 197 DF, p-value: < 2.2e-16

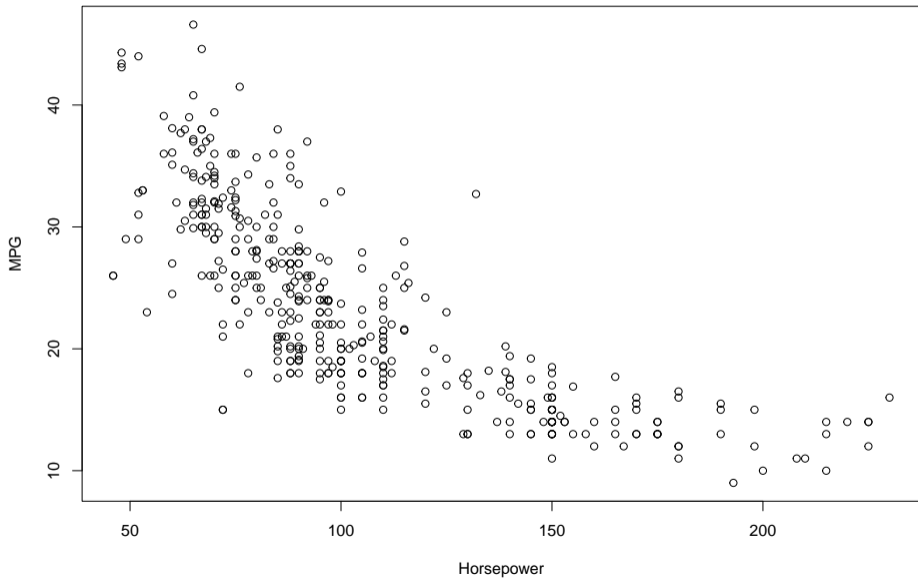
## Hierarchical Principal

In general, if we include an interaction term in a model, we also include the main effects *even if the  $p$ -values associated with the main effects are not significant.*

Consider Credit Balance predicted by Income and Student status.



- ▶ The interaction allows the model for students to have a different slope than the model for non-students, while the main effects model only allows for different intercepts.





# Nonlinear Relationships Between Predictors and Outcome

How can we deal with this using *linear* regression?

- ▶ The model fit requires the model to be linear *with respect to*  $\beta$ .
- ▶ This is much like including  $X_1X_2$  in the model by creating a “new variable” in the matrix  $X$ .
- ▶ Here, we just construct a “new variable”, say,  $X_1^2$  in  $X$ .

## Nonlinear Relationships Between Predictors and Outcome

```
mod5 <- lm(mpg ~ poly(horsepower, 2), data = Auto)
summary(mod5)
```

Coefficients:	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	23.4459	0.2209	106.13	<2e-16 ***
poly(horsepower, 2)1	-120.1377	4.3739	-27.47	<2e-16 ***
poly(horsepower, 2)2	44.0895	4.3739	10.08	<2e-16 ***

Residual standard error: 4.374 on 389 degrees of freedom

Multiple R-squared: 0.6876, Adjusted R-squared: 0.686

F-statistic: 428 on 2 and 389 DF, p-value: < 2.2e-16

