

TOPIC2: Multiple Linear Regression - Interaction Effects and Second-Order Models

© Thunida Ngankham 2022 modified by Paul Galperin

An Interaction Model with Quantitative Predictors

Consider the standard linear regression model with two variables,

additive model: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + e$

According to this model, if we increase X_1 , then Y will increase by β_1 . Notice that the presence of X_2 does not alter this statement; that is, regardless of the value of X_2 , a constant increase in X_1 will lead to a β_1 increase in Y . The assumption is also known as **homoscedasticity**: investigation of predictors assumes that the relationship between a given predictor variable and the response is independent of the other predictor variable.

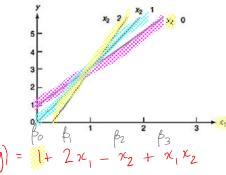
If the relationship between $E(Y)$ and X_1 depends on the values of the remaining X 's held fixed, then the first-order model is not appropriate for predicting Y . Interactions occur whenever two independent variables or a dependent variable is not constant over all of the values of the other independent variables. In particular, interactions arise when the effect of one variable depends on the value of another variable. Such a model includes the cross products of two or more X 's. Hence, the interaction model is

interaction model:

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

For example, suppose that the mean value $E(Y)$ of a response Y is related to two quantitative independent variables, X_1 and X_2 , by the model

$$E(Y) = 1 + 2X_1 - X_2 + X_1 X_2$$

A graph of the relationship between $E(Y)$ and X_1 for $X_2 = 0, 1$, and 2 is displayed in Figure 1.Note that the graph shows three nonparallel straight lines. You can verify that the slopes of the lines differ by substituting each of the values $X_2 = 0, 1$, and 2 into the equation.

$$E(Y) = 1 + 2X_1 - (1) + X_1(1) = 3X_1 (\text{slope} = 3)$$

$$E(Y) = 1 + 2X_1 - (2) + X_1(2) = -1 + 4X_1 (\text{slope} = 4)$$

Note that the slope of each line is represented by slope = $\beta_1 + \beta_3 x_2$. Thus, the effect on $E(Y)$ of a change in X_1 (i.e., the slope) now depends on the value of X_2 . When this situation occurs, we say that X_1 and X_2 interact. Note that the lines in Figure 1 are not parallel. The cross-product term, $X_1 X_2$, is called an interaction term, and the model $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + e$ is called an interaction model with two quantitative variables.

Testing for Interaction in Multiple Regression

For testing an interaction term in regression model, we use the **Partial Correlation Test** method.

- ① Set up H_0 & H_a
 - ② $d = 0.05$
 - ③ Calculate the test statistic & compare to t_{table} (df = p-1).
- where p is the total number of independent variables (including interaction terms).

Considering our advertising example, let's test the interaction term.

2

```

Advertising=read.table("Advertising.txt", header = TRUE, sep = "\t")
interactmodel<-lm(sale~tv+radio+radio*tv, data=Advertising)
summary(interactmodel)

## 
## Call:
## lm(formula = sale ~ tv + radio + radio * tv, data = Advertising)
## Residuals:
##   Min   1Q   Median   3Q   Max 
## -4.3566 -0.4028  0.1831  0.6948  1.5246 
## 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-02  12.699  <2e-16 *** 
## radio       2.886e-02  8.905e-03  3.241  0.0014***  
## radio*tv    1.086e-03  8.242e-05  1.327  0.1860    
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## 
## Residual standard error: 0.9436 on 196 degrees of freedom
## Multiple R-squared:  0.9678, Adjusted R-squared:  0.9673 
## F-statistic: 196.5 on 3 and 196 DF,  p-value: < 2.2e-16

glmmodel<-lm(sale~tv+radio, data=Advertising)
summary(glmmodel)

## 
## Call:
## lm(formula = sale ~ tv + radio, data = Advertising)
## Residuals:
##   Min   1Q   Median   3Q   Max 
## -6.3566 -0.4028  0.1831  0.6948  1.5246 
## 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-02  12.699  <2e-16 *** 
## radio       2.886e-02  8.905e-03  3.241  0.0014***  
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## 
## Residual standard error: 0.9436 on 196 degrees of freedom
## Multiple R-squared:  0.9678, Adjusted R-squared:  0.9673 
## F-statistic: 196.5 on 2 and 196 DF,  p-value: < 2.2e-16

```

Exactly the same - but simplified!

This is an additive model

They could be an oversimplification of the world. And we may wish to capture more complexity.

eg. It could be that x_1 & x_2 are related + therefore that we need to take this into account.

um

caveat

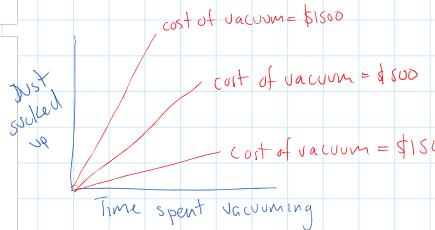
but not too related as this creates multicollinearity (we'll discuss this later)

In this course we're going to stop at 2-way interactions

→ 3 way + more are possible but too complicated to interpret

* We need to interpret our models + explain them to others. When they are too complex, it's hard!

INTERACTIVE CASE



$$E(Y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$

TERMS WITH x_1

$$= \beta_0 + (\beta_1 + \beta_3 x_2)x_1 + \beta_2 x_2$$

the slope becomes a function or the effect on x_1 is functionalcoeff. of x_1 of that interaction term

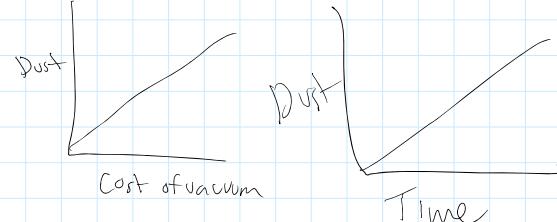
interaction term

term

term

term

ADITIVE



Why is it called multiple linear regression when sometimes we fit parameters that interact (i.e. multiply)?

 $\beta_0 + \beta_1 x_1 + \beta_2 x_2$ } linear combination

the model is

linear in the parameters

aka → The model assumes the relationship between $y + x_1$ is linear

* Although interaction terms seems non-linear it is still consider a linear model.

both additive interactive

Different R formula interfaces

$$E(Y) = \beta_0 + \beta_1(x) + \beta_2(y) + \beta_3(x)(y)$$

this can be used to understand

$$\textcircled{1} \quad y = a + b + a:b \quad \xrightarrow{\text{try to understand}}$$

```

## Call:
## lm(formula = sales ~ tv + radio, data = Advertising)
## 
## Residuals:
##   Min   1Q Median   3Q   Max 
## -6.3566 -0.4028  0.1831  0.5948  1.5246 
## 
## 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 8.958e-03 3.241 0.0014 **  
## tv:radio    1.098e-03 8.294e-08 130.727 <2e-16 ***  
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## 
## Residual standard error: 0.9438 on 196 degrees of freedom
## Multiple R-squared:  0.9675, Adjusted R-squared:  0.9673 
## F-statistic: 196.3 on 3 and 196 DF, p-value: < 2.2e-16
#System.out.println("intercodel2=lm(sale-(tv+radio)^2, data=Advertising)
summary(intercodel2)
```

```

## Call:
## lm(formula = sales ~ (tv * radio)^2, data = Advertising)
## 
## Residuals:
##   Min   1Q Median   3Q   Max 
## -6.3566 -0.4028  0.1831  0.5948  1.5246 
## 
## 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 8.958e-03 3.241 0.0014 **  
## tv:radio    1.098e-03 8.294e-08 130.727 <2e-16 ***  
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## 
## Residual standard error: 0.9439 on 196 degrees of freedom
## Multiple R-squared:  0.9675, Adjusted R-squared:  0.9673 
## F-statistic: 196.3 on 3 and 196 DF, p-value: < 2.2e-16

```

As you can see from the output, $t_{\text{radio}} = 20.727$ with the $p\text{-value} < 0.0001$, indicating that we should clearly reject the null hypothesis which means that we should definitely add the interaction term to the model at $\alpha = 0.05$. Model $2 = \text{Intercodel2}$ is better than model 1 because the response variable is explained by the interaction model compared to only 5.8562 for the additive model that predicts sales using TV and radio without an interaction term.

Note that from the additive model assume that the effect on sales of TV advertising is independent of the effect of radio advertising. This assumption might not be true. For example, spending money on TV advertising may increase the effectiveness of radio advertising on sales. In [MLRPage104.html](#), [MLRPage105.html](#) it is shown that ~~interaction effects~~ in statistics it is referred to as ~~interaction effects~~

Interpreting Coefficients of Predictor Variables

Notice that the model also can be rewritten as

$$Y = \beta_0 + (\beta_1 X_1 + \beta_2 X_2 + \epsilon) \quad \text{Rewrite by factoring to know the effect in } X_1$$

where $= \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$

$$Y = \beta_0 + \beta_1 X_1 + (\beta_2 + \beta_1 X_1) X_2 + \epsilon \quad \text{Rewrite by factoring to know effect on } X_2$$

where $= \beta_0 + \beta_2 X_2$

For economic purposes, we are interested in the coefficient of a factor, we want to predict the number of workers produced based on the number of production lines and the total number of workers. It seems likely that the effect of increasing the number of production lines will depend on the number of workers, since if no workers are available to operate the lines, then increasing the number of lines will not

4

relying for an interaction term from "int principle" or object knowledge

increase production. This suggests that it would be appropriate to include an interaction term between lines and workers in a linear model to predict units. Suppose that when we fit the model, we get

$$\hat{units} = 1.9 + 3.4X_1 + 0.7X_2 + 1.4(X_1 \cdot X_2) \quad \text{makeup fit of data}$$

In other words, adding an additional worker will produce 3.4 more units, and an additional line produced by 2.4 + 1.4 workers. Hence, the more workers we have, the stronger will be the effect of lines. Let's return to the Advertising example. A linear model that uses radio, TV, and an interaction between the two to predict sales takes the form

$$Sale = \beta_0 + \beta_1 TV + \beta_2 Radio + \beta_3 (TV \cdot radio) + \epsilon \quad \text{to find effect on TV}$$

We can interpret the coefficient β_3 as spending additional 1,000 dollars on TV advertising leads to an increase in sales by approximately $\beta_3 + \beta_2 \text{radio}$ units.

The results of the output strongly suggest that the model that includes the interaction term to the model that excludes the interaction term effects

The coefficient estimates in the model suggest that an increase in TV advertising of 1,000 dollars is associated with an increase in sales of $(\beta_3 + \beta_2 \text{TV}) \times 1,000 = 29 + 117V$ units.

In this example, the p -values associated with TV, radio, and the interaction term are all statistically significant and so it is obvious that all three variables should be included in the model. However, it is sometimes the case that an interaction term has a very small p -value but the associated main effects (in this case, TV and radio) do not.

the reason of an interaction term in a model is to show that the effect of one predictor variable depends on the value of another predictor variable

Caution That is, if the interaction between X_1 and X_2 seems important, then we should include both X_1 and X_2 in the model even if their coefficient estimates have large p -values.

The rationale for this principle is that if X_1, X_2, \dots, X_k are the variables in the model, then whether or not the coefficients of X_1, X_2, \dots, X_k are statistically zero is of little interest. X_1, X_2, \dots, X_k are typically correlated with X_1, X_2, \dots, X_k so leaving them out tends to favor the selection of the intercept.

In-class Practice Problem 4

From the condominium problem, do the data provide sufficient evidence to indicate that the interaction term need to be added in the model? If you had to compare additive models with the interaction model, which model would you choose? Explain

In-class Practice Problem 5

Data on last year's sale (Y in 100,000s dollars) in 40 sales districts are given in the sales.csv file. This file also contains

promotional expenditures (X_1 ; in 1,000s dollars), the number of active accounts (X_2), the number of competing banks (X_3) and the district potential (X_4 , coded) for each of the district (OMIT THIS VARIABLE FOR NOW).

- Find the best fit additive to predict sales using some or all of the variables X_1, X_2, X_3 only.
- Find the best fit model with interaction terms (if needed) using some or all of the variables X_1, X_2, X_3 .
- Which model would you choose? Explain.

5

- Once you obtain the best fit model, interpret the regression coefficient for X_3 (Hint: it will interact with another variable).

Multiple Regression with Qualitative (Dummy) Variable Models

Multiple regression models can also be written to include qualitative (or categorical) independent variables. Qualitative variables, unlike quantitative variables, cannot be measured on a numerical scale. Therefore, we must code the values of the qualitative variable (called levels) as numbers before we can fit in the model. These coded qualitative variables are called **dummy variables** or **categorical predictor variables**, since the number assigned to the various levels are arbitrarily selected.

Dummy Coding:
Because categorical predictor variables cannot be entered directly into a regression model and be meaningfully interpreted, some other methods of dealing with information of this type must be developed. In general, a categorical variable with k levels will be transformed into $k-1$ variables each with two levels. For example, if a categorical variable had k levels, then $k-1$ new variables could be constructed that would contain

$$E(y) = \beta_0 + \beta_1(x_1) + \beta_2(x_2) + \beta_3(x_3)$$

that can be written down

$$\textcircled{1} \quad y = a + b + a:b \quad \text{--- expect to understand}$$

AVOID because this will give you unexpected results with 3 or more levels

$$\textcircled{2} \quad y = a + b + b^2 \quad \text{--- correct approach when you have a lot of variables}$$

3 or 2 way interactions

$$\textcircled{3} \quad y = a + b + c + d + e + f + g + h + i + j + k + l + m + n + o + p + q + r + s + t + u + v + w + x + y + z$$

gives you each variable alone plus pairwise interactions.

R side

PMS
PMS
PMS

key idea
Not just $(TV) + (Radio)$
Rather it's $(TV) + (Radio) + [How (Radio) affects (TV)]$
[How (TV) affects (Radio)]

key synergies for interaction

Some changes with X_3 , the effect of X_3 on Y is no longer constant, depending X_3 will change the effect of X_1 on Y. The model also can be rewritten as

$$Y = \beta_0 + (\beta_1 X_1 + \beta_2 X_2 + \epsilon) \quad \text{Functional co efficient}$$

For economic purposes, we are interested in the coefficient of a factor, we want to predict the number of workers produced based on the number of production lines and the total number of workers.

It seems likely that the effect of increasing the number of production lines will depend on the number of workers, since if no workers are available to operate the lines, then increasing the number of lines will not

$y = A + B + A:B$
WITH ADDING
 A — — — $p = 0.001$
 B — — — $p = 0.002$

WITH INTERACTION MODEL

By hierarchical principle we must keep B in model because A is a signifiant A:B — $p = 0.01$

$A:B$ — $p = 0.01$

A — — — $p = 0.025$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.001$

B — — — $p = 0.002$

$A:B$ — $p = 0.01$

A — — — $p = 0.00$

Qualitative variables, unlike quantitative variables, cannot be measured on a numerical scale. Therefore, we must code the qualitative variable (called levels) as numbers before we can fit in the model. These coded qualitative variables are called **dummy variables** or **categorical predictor variables**, since the number assigned to the variable levels are arbitrarily selected.

Dummy Coding

Binary predictor variables can be entered directly into a regression model and be meaningfully interpreted. Some other methods of dealing with information of this type must be developed. In general, a categorical variable with k levels will be transformed into $k - 1$ variables each with two levels. For example, if a categorical variable had six levels, then five dichotomous variables could be constructed that would contain the same information as the single categorical variable. The process of creating dichotomous variables from categorical variables is called **dummy coding**.

Dummy Coding: Coding with two levels to the case of dummy coding is when a categorical variable has two levels, assigning zero and one to the variable.

For example, the Credit data set records balance (average credit card debt for a number of individuals) as well as several quantitative predictors: age, cards (number of credit cards), education (years of education), income (in thousands of dollars), limit (credit limit), and rating (credit rating). In addition to these quantitative variables, there are two qualitative variables: gender (Male/Female) and marital status (Married/Not Married), and ethnicity (Caucasian/Asian/Black/American Indian). Data are provided in **Credit.csv**. Suppose that we wish to investigate differences in credit card balance between males and females. Based on the gender variable, we can create a dummy variable with 0 as male and 1 as female.

```
credit<-read.csv('credit.csv',header = TRUE)
#  
#  
## number Income Limit Rating Cards Age Education Marital_Status  
## | 1 14.891 3606 283 2 34 11 Native No Married  
## | 2 106.029 6645 483 3 82 15 Female Yes Yes  
## | 3 104.593 7075 514 4 71 11 Native No No  
## | 4 49.125 3025 483 5 26 16 Native No Yes  
## | 5 55.889 4897 357 2 68 16 Native Yes Yes  
## | 6 80.140 3027 569 4 77 10 Native No No  
## | 7 100.000 3027 569 4 77 10 Native No No  
## | 8 Caucasian  
## | 9 Asian  
## | 10 Black  
## | 11 American Indian  
## | 12 Native  
## | 13 Caucasian  
## | 14 Asian  
## | 15 Black  
## | 16 American Indian  
## | 17 Native  
dummymodel1<-lm(Balance ~ factor(Gender), data=credit)
summary(dummymodel1)
```

"**factor**" means a categorical variable. forces R to handle it as a dummy variable

Residuals:

```
## Min 1Q Median 3Q Max
## -829.54 -458.38 -60.17 334.71 1489.20
## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) 33.13 15.389 <2e-16 ***
## factor(Gender)Female 46.05 0.429 108.000 (not significant) if we were being pedantic we should interpret this as not significant
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '
## Residual standard error: 460.2 on 398 degrees of freedom
## Multiple R-squared: 0.000461, Adjusted R-squared: -0.00205
## F-statistic: 0.1838 on 1 and 398 DF, p-value: 0.6685
optical<-lm(Balance ~ factor(Gender), data=credit)
summary(optical)
```

Rfunction

factor() command will make sure that R knows that your variable is categorical. This is especially useful if your categories are indicated by integers, otherwise **lm()** might interpret the variable as continuous.

$$\widehat{balance}_{ij} = \begin{cases} \beta_0 + \beta_1 + \epsilon & \text{if } i^{\text{th}} \text{ person is female} \\ \beta_0 + \epsilon & \text{if } i^{\text{th}} \text{ person is male} \end{cases}$$

Interpreting Coefficients of Predictor Variables

Can be interpreted as the average credit card balance among females $\beta_1 = \beta_0(a)$ what is difference

Can be interpreted as the average credit card balance among males $\beta_1 = \beta_0(b)$

Let's fill in the values from R output above

$$\widehat{Y}_i = 509.80 + 19.73X_{1i}$$

$$\widehat{balance}_{ij} = \begin{cases} 509.80 + 19.73 & \text{if } i^{\text{th}} \text{ person is female} \\ 509.80 & \text{if } i^{\text{th}} \text{ person is male} \end{cases}$$

From the output, the coefficient estimates and other information associated with the model are provided.

The coefficient estimate for gender is 19.73, which indicates that females are estimated to carry 19.73 additional debt for a total of $(509.80 + 19.73) = 529.53$ dollars. However, we notice that the p-value for the dummy variable is very high. This indicates that there is no statistical evidence of a difference in average credit card balance between the genders.

Inclass Practice Problem 6

Suppose that we wish to investigate differences in credit card balance between marital status. Based on the Marital variable, we can create a dummy variable which is NO if 1 is Yes. Create a simple linear regression model to predict the credit card balance by using the gender variable. What are the regression coefficients for this model? How do you interpret the regression coefficient? Ignore the individual t-test output.

Dummy Coding with three levels: When the categorical variable has three levels, it is converted to two dichotomous (dummy) variables.

For example, there is always a certain curiosity and controversy surrounding professor salaries and whether the ever-increase not paid enough. A university would like to study the effects of ranks and department on salary. The observations were randomly chosen from different departments. The data are provided in **salary.csv**.

gender: {Male, Female}

rank: {Assistant, Associate, Full}

Dept: Department {B-Family, B-Biology, B-Business}

year-Years since Merit Rank

merit-Average Merit Ranking

The variable **Dept** has three levels: Family, Biology, and Business. Variable **Dept** could be dummy coded into two variables, one called Biology and one called Business. The variable **rank** will also have three levels and will be coded into two dummy variables: Assistant Prof and Full Prof. The dummy coding is represented below.

Before considering both predictors, let practice how to interpret the regression coefficients for each categorical variable.

For example, considering only rank variable with three levels

```
salary<-read.csv('salary.csv',header = TRUE)
head(salary)
```

```
## salary gender rank Dept year merit
## 1 38 0 B-Assistant 2 4.38
## 2 58 1 B-Assistant 3 4.38
## 3 60 0 B-Assistant 4 4.38
## 4 30 1 B-Assistant 5 2.54
## 5 50 1 B-Assistant 6 2.06
```

Qualitative / Categorical variables

	x_1	x_2
- EXCERCISE	0	0
- AVERAGE	1	0
- SUB-SUPERIOR	0	1

3 levels $\rightarrow k=3$ DUMMIES VARIABLES

dummys vars $k-1$ columns

K=2 2-1 dummy variables

MALE 0

FEMALE 1

$$y_i = \beta_0 + \beta_1 \text{gender}_i + \epsilon_i$$

this model asks simply is there an effect of gender on the response (y)

Male 0 1
Female 1 0
Non-binary 0 1

k=3 : 3-1 = 2 dummy variables			
Dummy Coding with three levels:			
dummy variable:	Department	Biology	Business
Business		1	0
Biology		0	1
Family Studies		0	0

dummy variable:	Rank	Associate Prof	Full Prof
Associate Prof		0	0
Full Prof		1	0
Assistant Prof		0	1

How does R know which is the 0 0 case (reference)?
levels (factor salary/rank)

Figure 1: Dummy Coding with three levels

```
dummymodel<-lm(salary~factor(rank),data=salary)
summary(dummymodel)

## Call:
## lm(formula = salary ~ factor(rank), data = salary)
## Residuals:
##   Min   1Q   Median   3Q   Max 
## -18.875 -5.799  0.000  5.353 23.125 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept)  2.259  18.686  0 <2e-16 ***
## factor(rank2) 4.005  2.640  1.518  0.1265    
## factor(rank3) 3.830  3.884  0.000600 *** 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## Residual standard error: 8.749 on 27 degrees of freedom
## Multiple R-squared:  0.3881, Adjusted R-squared:  0.3428 
## F-statistic: 8.963 on 2 and 27 DF,  p-value: 0.001319
```

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

```
oops! we didn't specify factor! →
(+) treat rank as a numerical variable + not a category
NOR WHAT WE WANT.

## lm(formula = salary ~ factor(rank), data = salary)
## Residuals:
##   Min   1Q   Median   3Q   Max 
## -19.9017 -5.2168  0.1139  5.8639 22.0983 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 34.871  3.694  9.466 3.19e-10 ***
## rank        0.677  1.030  4.080 0.000339 *** 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## Residual standard error: 8.697 on 26 degrees of freedom
## Multiple R-squared:  0.3726, Adjusted R-squared:  0.3698 
## F-statistic: 16.65 on 1 and 26 DF,  p-value: 0.0003389
```

$$\text{rank refers to type of software}$$

$$\text{salary}_i = \begin{cases} \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon & \text{if } i^{\text{th}} \text{ person is ranked as Associate Prof} \\ \beta_0 + \beta_1 + \beta_2 + \epsilon & \text{if } i^{\text{th}} \text{ person is ranked as Full Prof} \\ \beta_0 + \epsilon & \text{if } i^{\text{th}} \text{ person is ranked as Assistant Prof} \end{cases}$$

Good model based on rank

individual t-test

$$\text{salary}_i = \begin{cases} \beta_0 & ① \\ \beta_0 + \beta_1 & ② \\ \beta_0 + \beta_1 + \beta_2 & ③ \end{cases}$$

if Assist (sub in 0, 0)
if Associate (sub in 1, 0)
if Full (sub in 0, 1)

What should I do now if one category is not significant?

A **
B [] no sign!

You need to leave them all in (all levels of factor)

an important point

can be interpreted as the average salary for Assistant Professor position.
 β_1 as the difference in average salary between Associate Professor and Assistant Professor.
 β_2 as the difference in average salary between Full Professor and Assistant Professor.

10

can be interpreted as the average salary for Full Professor position.

In-class Practice Problem 7

There is always a certain curiosity and controversy surrounding professors' salaries and whether they are overpaid or not paid enough. A university would like to study the effects of ranks and departments on salaries. 30 observations were randomly chosen from 3 different departments. The data are provided in the salary.csv data file. The columns in the data file are:

Example: There is always a certain curiosity and controversy surrounding professors' salaries and whether they are overpaid or not paid enough. A university would like to study the effects of ranks and dept on the salaries (thousands of dollars). 30 observations were randomly chosen from 3 different departments. The data are provided in salary.csv data file.

```
salary.read.csv("salary.csv",header = TRUE)
dummymodel<-lm(salary ~ factor(rank) + factor(dept), data = salary)
summary(dummymodel)

## Call:
## lm(formula = salary ~ factor(rank) + factor(dept), data = salary)
## Residuals:
##   Min   1Q   Median   3Q   Max 
## -11.243 -3.333 -0.049  2.350 20.256 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 34.939  2.308 14.878 0.6204 ***
## factor(rank2) 2.029  2.477  0.82165 *** 
## factor(rank3) 15.194  2.797  5.453 1.22e-05 ***
## factor(dept2) 10.502  2.972  3.533 0.001624 ** 
## factor(dept3) 13.381  2.752  4.881 4.48e-08 *** 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## Residual standard error: 6.368 on 25 degrees of freedom
## Multiple R-squared:  0.6995, Adjusted R-squared:  0.6518 
## F-statistic: 14.57 on 4 and 25 DF,  p-value: 2.85e-06
```

factor() command will make sure that R knows that your variable is categorical. This is especially useful if your categories are indicated by integers, otherwise function *lm()* will interpret the variable as continuous.

$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \beta_4 X_{i4} + \epsilon$

two factors or categorical variables in the same model.

factor(): command will make sure that R knows that your variable is categorical. This is especially useful if your categories are indicated by integers, otherwise function lm() will interpret the variable as continuous.

$$Y_1 = \beta_0 + \beta_1 X_{11} + \beta_2 X_{12} + \beta_3 X_{13} + \beta_4 X_{14} + \epsilon$$

$$Y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \epsilon$$

A graph like this or the one above
is called a scatter plot.
The points represent individual observations.

Interpreting Coefficients of Predictor Variables

β_0 can be interpreted as the average salary of an Assistant Prof from Family Studies dept.

β_1 can be interpreted as the average difference in salary between an Associate Prof (hold dept).

β_2 can be interpreted as the average difference in salary between an Assistant Prof and a Full Prof (hold dept).

.

Interaction Effect in Multiple Regression with both Quantitative and Qualitative (Dummy) Variable Models

In previous topics, we considered Multiple Regression models for both quantitative and qualitative variables. We also discussed an interaction in Multiple Regression for quantitative variables. However, the concept of interaction goes just as well to qualitative variables, or a combination of quantitative and qualitative variables. In fact, an interaction between a quantitative variable and a qualitative variable has a particularly nice interpretation.

Consider the example and suppose that we wish to predict balance using the associated and student variables. In the absence of an interaction term, the model takes the form

$$Y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

balance_i = $\beta_0 + \beta_1 \text{Income}_i + \beta_2$ if i^{th} person is a student
balance_i = $\beta_0 + \beta_2$ if i^{th} person is not a student

fit thus model into an interaction

balance_i = $\beta_0 + \beta_1 \text{Income}_i + \beta_2$ if i^{th} person is a student
balance_i = $\beta_0 + \beta_2$ if i^{th} person is not a student

student

no 0
yes 1

```
credit=read.csv("credit.csv",header = TRUE)
mixmodel<-lm(balance~Income+factor(student), data=credit)
summary(mixmodel)
```

```
##
```

```
## Call:
```

```
## lm(formula = Balance ~ Income + factor(Student), data = credit)
```

```
## Residuals:
##   Min   1Q   Median   3Q   Max 
## -762.37 -331.38 -45.04 323.60 818.28 
## 
## Coefficients:
##   Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 32.4672   6.5055  2.34e-10 ***
## income      5.9843   0.8568  10.751 <2e-16 ***
## factor(Student)Yes 65.3168   5.869  9.79e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 391.8 on 397 degrees of freedom
## Multiple R-squared:  0.2775, Adjusted R-squared:  0.2738 
## F-statistic: 76.22 on 2 and 397 DF, p-value: < 2.2e-16
```

balance_i = $\beta_0 + \beta_1 \text{Income}_i + \beta_2$ if i^{th} person is a student
balance_i = $\beta_0 + \beta_2$ if i^{th} person is not a student

student

no 0
yes 1

Fit this no interaction model

Substitute in parameters from Root of

Notice that this amounts to fitting two parallel lines to the data, one for students and one for non-students. The lines for students and non-students have different intercepts, $\beta_0 + \beta_2$ versus β_0 , but the same slope, β_1 . This is illustrated in the below. The two parallel lines have the same slope, β_1 , but different intercepts, $\beta_0 + \beta_2$ and β_0 . A balance of a one-unit increase in income increases by β_1 for both groups, but the increase is larger for students than for non-students. Depending on whether you're a student or not

```
library(ggplot2)
credit=read.csv("credit.csv",header = TRUE)
mixmodel<-lm(balance~Income+factor(Student), data=credit)
#For student y=391.8124+5.9843*Income
#For nonstudent y=65.3168+5.9843*Income
nonstudent=function(x){coef(mixmodel)[2]+x*coef(mixmodel)[1]}
student=function(x){coef(mixmodel)[2]+x*coef(mixmodel)[1]+coef(mixmodel)[0]}
ggplot(credit,aes(x=Income,y=balance))+
  stat_function(fun=nonstudent,geom="line",color=scales::hue_pal(2)[2])
  stat_function(fun=student,geom="line",color=scales::hue_pal(2)[1])
```

RANK	x ₁	x ₂	DEPT	x ₃	x ₄
ASSIST	0	0	F.S.	0	0
ASOC	1	0	BIO	1	0
FULL	0	1	BUSINESS	0	1

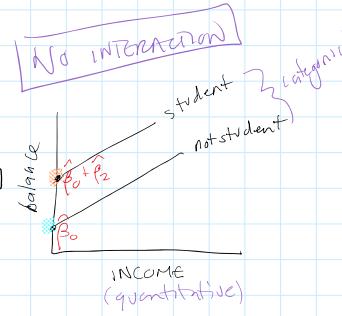
All the ways of getting β
 ③ → ④
 ⑤ → ④
 ① → ④

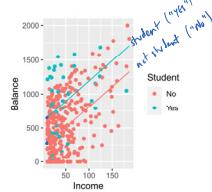
Now we have qualitative (categorical) and quantitative variables in the same model + we are interacting them.

Interaction qual × qual
 (qual × quant)
 Quant × quant — done already

If STUDENT balance = $(\beta_0 + \beta_2) + \hat{\beta}_1 [\text{income}]$
 Not STUDENT balance = $(\beta_0) + \hat{\beta}_1 [\text{income}]$

$$y = \text{int} + \text{slope} \cdot x$$





This limitation can be addressed by adding an interaction variable, created by multiplying income with the dummy variable for student. Our model now becomes

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i1} X_{i2} + \epsilon_i$$

interactions are created by multiplying

$$\text{balance}_i = \beta_0 + \beta_1 x_{i1} \underset{(1)}{\text{if } i^{\text{th}} \text{ person is a student}} + \beta_2 x_{i2} \underset{(0)}{\text{if } i^{\text{th}} \text{ person is not a student}}$$

show me the algebra, pal!

$$\text{Assuming a student} \therefore x_{i2} = 1 \text{ and } x_{i1} = [\text{income}]$$

$$\Rightarrow y_i = f_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i1} x_{i2}$$

$$= f_0 + \beta_1 [\text{income}] + \beta_2 (1) + \beta_3 [\text{income}] (1)$$

$$= (\beta_0 + \beta_2) + (\beta_1 + \beta_3) [\text{income}]$$

int slope

interaction between quant qual

fit the model

14

write out the model statement

please practice this on your own (i.e. writing submodel levels)

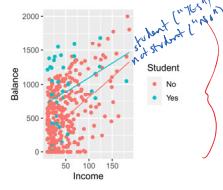
$$Y_i = 200.6232 + 6.2182 X_{i1} + 476.6758 X_{i2} - 1.9992 X_{i1} X_{i2} + \epsilon_i$$

$$\text{balance}_i = 200.6232 + 6.2182 \text{income}_i + \begin{cases} 476.6758 - 1.9992 \text{income}_i & \text{if } i^{\text{th}} \text{ person is a student} \\ 0 & \text{if } i^{\text{th}} \text{ person is not a student} \end{cases}$$

$$\widehat{\text{balance}}_i = \begin{cases} 200.6232 + 6.2182 \text{income}_i & \text{if } i^{\text{th}} \text{ person is a student} \\ 200.6232 + (6.2182 - 1.9992) \text{income}_i & \text{if } i^{\text{th}} \text{ person is not a student} \end{cases}$$

$$\widehat{\text{balance}}_i = \begin{cases} 200.6232 + 4.2191 \text{income}_i & \text{if } i^{\text{th}} \text{ person is a student} \\ 200.6232 + 6.2182 \text{income}_i & \text{if } i^{\text{th}} \text{ person is not a student} \end{cases}$$

```
library(ggplot2)
credit<-read.csv("credit.csv",header = TRUE)
mixmodel<- lm(Balance~Income+factor(Student),data=credit)
summary(mixmodel)
#or nonstudent y=200.6232+6.2182*income
student=function(i)(4.2191*677.2992)
nonstudent=function(i)(6.2182*677.2992)
nonstudent=function(i)(coef(mixmodel)[2]*coef(mixmodel)[1])
ggplot(credit, aes(x=Income, y=Balance, color=factor(Student))) + geom_point() +
stat_function(fun=student, geom="line",color=scales::husl_pal(2)[1]) +
stat_function(fun=nonstudent, geom="line",color=scales::husl_pal(2)[2])
```



15

Discuss the results for the interaction term. we have two different regression lines for the students and the non-students. Both the two regression lines have different intercepts, $\beta_0 + \beta_1$ versus β_0 , as well as different slopes, $\beta_1 + \beta_2$ versus β_1 . This allows for the possibility that changes in income may affect the credit card balances of students and non-students differently. The output shows the estimated relationships between income and balance for students and non-students in the model. We note that the slope for students (4.219) is lower than the slope for non-students (6.218). This suggests that increases in income are associated with smaller increases in credit card balance among students as compared to non-students.

Inclass Practice Problem 8

From the credit card example, use the lm() function to perform the best fit model. How would you interpret the regression coefficient (if possible)? Would you recommend this model for predictive purpose?

For the inclass practice problem, you can see that it is quite complicated (possible) to interpret regression coefficients as they are not all predictors in the model. However, let's practice the example below.

```
credit=read.csv('credit.csv', header=TRUE)
credit$balance=balance*factor(credit$income)
credit$student=factor(credit$student)
credit$rating=factor(credit$rating)
summary(credit)
summary(mimmo$ch42)
```

##

Call:

lm(formula = Balance ~ Rating * Income + factor(Student) + factor(Student) *

Rating + factor(Student) * Income + Rating * Income, data = credit)

Residuals:

Min 1Q Median 3Q Max

-207.799 -74.852 -3.999 74.055 253.340

Coefficients:

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

(Intercept) 4.712e+00 2.053e+01 -22.485 < 2e-16 ***

Rating 3.701e+00 6.312e-02 58.642 < 2e-16 ***

Income -0.925e+00 4.696e-01 -21.136 < 2e-16 ***

factor(Student)Yes 1.862e+02 4.442e+01 4.193 3.41e-05 ***

factor(Rating)1 1.862e+02 4.442e+01 4.193 3.41e-05 ***

Income:factor(Student)Yes -0.275e+00 6.738e-01 -4.121 4.59e-05 ***

Rating:Income 4.140e-03 7.098e-04 5.832 1.14e-08 ***

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

##

Residual standard error: 96.05 on 393 degrees of freedom

Multiple R-squared: 0.857, Adjusted R-squared: 0.8664

F-statistic: 1450 on 6 and 393 DF, p-value: < 2.2e-16

balance = $\beta_0 + \beta_1 \text{Rating} + \beta_2 \text{Income} + \beta_3 \text{Student}$
 $+ \beta_4 \text{Rating} * \text{Student} + \beta_5 \text{Income} * \text{Student} + \beta_6 \text{Rating} * \text{Income} + \epsilon$

if a person is a student
 $balance = \beta_0 + \beta_1 \text{Rating} + \beta_2 \text{Income} + \beta_3(1) + \beta_4 \text{Rating} * (1)$
 $+ \beta_5 \text{Income} * (1) + \beta_6 \text{Rating} * \text{Income} + \epsilon$

16

if a person is not a student
 $balance = \beta_0 + \beta_1 \text{Rating} + \beta_2 \text{Income} + \epsilon$

ALSO KNOWN AS
 POLYNOMIAL REGRESSION!

A Quadratic (Second Order) Model with Quantitative Predictors

All of the models discussed in the previous sections proposed straight-line relationships between $E(y)$ and each of the independent variables in the model. In this section, we consider a model that allows for curvature in the relationship. This model is a second-order model because it will include an X^2 term. Here, we consider a model that includes only one independent variable X . The form of this model, called the quadratic model, is

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \epsilon$$

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_1^2$$

The term involving X_1^2 , called a quadratic term (or second-order term), enables us to hypothesize curvature in the graph of the model relating Y to X_1 . Graphs of the quadratic model for two different values of β_2 are shown in the figure below. When the curve opens upward, the sign of β_2 is positive (see Figure 2 (a)); when the curve opens downward, the sign of β_2 is negative (see Figure 2 (b)).

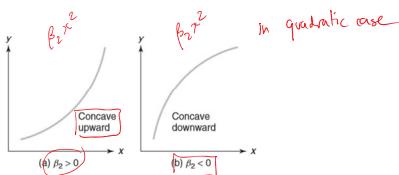
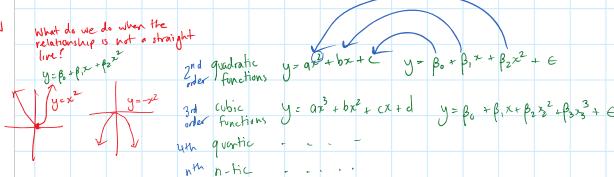


Figure 2: Graphs for two quadratic models

Interpretation of the regression coefficients

The interpretation of the estimated coefficients in a quadratic model must be taken cautiously. $\hat{\beta}_1$ can be meaningfully interpreted only if the range of the independent variable includes zero; that is, if $X_1 = 0$ is included in the sampled range of X_1 . $\hat{\beta}_1$ no longer represents a slope in the presence of the quadratic term X_1^2 . The estimated coefficient of the first-order term X_1 will not, in general, have a meaningful interpretation in the quadratic model.

17



How do higher order models fit into our schema of modelling?

- ① build an additive model
- ② check for interactions
- ③ Ask: can we improve the fit by adding higher-order terms?

How DO WE KNOW WHICH POWER TO USE? (quadratic, cubic etc?)

*Start at linear

and increase the degree of the polynomial until we find the highest degree is NOT SIGNIFICANT

CAVEAT EMPTOR

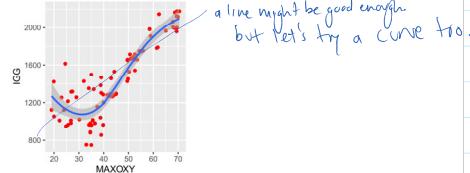
We lose the ability to easily interpret the β coefficients

β_2 , the sign of the coefficients, β_2 , of the quadratic term, X_2^2 , is the indicator of whether the curve is concave downward (U-shaped) or concave upward (bow-shaped). A negative β_2 implies downward concavity, while in this example (Figure 2), a positive β_2 implies upward concavity. Rather than interpreting the numerical value of β_2 itself.

Example A physiologist wants to investigate the impact of exercise on the human immune system. The physiologist theorizes that the amount of immunoglobulin Y in blood (called IgG, an indicator of long-term immunity, milligrams) is related to the maximal oxygen uptake x (a measure of aerobic fitness level, milliliters per kilogram). The data file is provided in **AEROBIC.CSV**. The **geom_smooth()** function in ggplot2 is the best model to fit the data.

```
library(ggplot2) #using ggplot2 for data visualization
aerobicdata=read.csv("AEROBIC.csv", header = TRUE)
ggplot(data=aerobicdata,aes(y=IgG,x=(MAXOXY))) + geom_point(color="red") +
  geom_smooth() # it's not "w" -> this means "just smooth it"
```

'geom_smooth()' using method = "loess" and formula 'y ~ x'



```
simplemodel=lm(IgG~MAXOXY,data=aerobicdata)
```

```
## 
## Call:
## lm(formula = IgG ~ MAXOXY, data = aerobicdata)
## 
## Residuals:
##   Min    1Q Median    3Q   Max 
## -479.11 117.30  58.04 114.38 426.34 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 388.854   69.561   5.735 1.38e-07 ***
## MAXOXY      1.668  16.120  0.10364    
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 201.4 on 87 degrees of freedom
```

18

READ THIS

R (statistic ASIDS)

- $(A+B)^2 \Rightarrow$ main effects + interactions
- $I(A^2) \Rightarrow$ Treat A as a second order term
- ↓ enters
- A^2 into the model

Do NOT DO $A \times A$

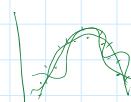
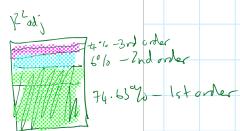
as long as the highest order term is significant we keep lower order terms even if they are not. hierarchical principle.

```
## Multiple R-squared:  0.7492, Adjusted R-squared:  0.7365 
## F-statistic: 250.6 on 1 DF, p-value: < 2.2e-16
quadmodel=lm(IgG~MAXOXY+I(MAXOXY^2), data=aerobicdata)
summary(quadmodel) # now we specify a second-order quadratic item
```

```
## 
## Call:
## lm(formula = IgG ~ MAXOXY + I(MAXOXY^2), data = aerobicdata)
## 
## Residuals:
##   Min    1Q Median    3Q   Max 
## -439.91 -86.43 -30.15 139.15 517.61 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 1270.41137 186.19900  6.823 1.18e-09 ***
## I(MAXOXY^2) -18.67954   8.52630  -2.158 0.03864 *  
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 176.6 on 86 degrees of freedom
## Multiple R-squared:  0.805, Adjusted R-squared:  0.8004 
## F-statistic: 177.5 on 2 and 86 DF, p-value: < 2.2e-16
cubemodel=lm(IgG~MAXOXY+I(MAXOXY^2)+I(MAXOXY^3), data=aerobicdata)
summary(cubemodel)
```

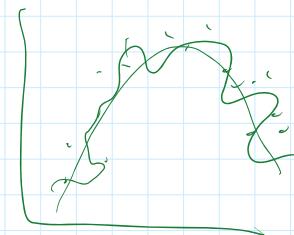
```
## 
## Call:
## lm(formula = IgG ~ MAXOXY + I(MAXOXY^2) + I(MAXOXY^3), data = aerobicdata)
## 
## Residuals:
##   Min    1Q Median    3Q   Max 
## -366.7 -100.1 -12.5 103.6 496.1 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 3.505e+03 5.015e+02  6.982 6.03e-10 ***
## I(MAXOXY) -11.964e-01 3.125e-01 -3.777 0.000117 *** 
## I(MAXOXY^2) -2.209e-02 6.379e-03 -3.372 0.000753 *** 
## I(MAXOXY^3) 1.860e-03 5.959e-04 3.110 0.001861 **  
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 169.9 on 85 degrees of freedom
## Multiple R-squared:  0.844, Adjusted R-squared:  0.84 
## F-statistic: 188.3 on 3 and 85 DF, p-value: < 2.2e-16
forthmodel=lm(IgG~MAXOXY+I(MAXOXY^2)+I(MAXOXY^3)+I(MAXOXY^4), data=aerobicdata)
summary(forthmodel) # should stay at forthmodel because all variables are not significant.
```

19



risk of overfitting

IF you increase order and R_{adj}^2 does not increase meaningfully go for the simpler option (i.e. don't increase to next order)



as our order increased there may be a tendency to OVERFIT

LOSS OF GENERALITY OF MODEL.

```

##  data = aerobicdata)
##  Residuals:
##    Min      1Q  Median      3Q     Max 
## -362.89 -104.07   -8.92   98.60  481.75 
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 2.012e+03 1.586e-03 1.261   0.211    
## MAFNTR      4.587e-02 1.586e-03 2.907   0.037    
## I(MAFNTR^2) -1.258e+00 5.947e+00 -0.211   0.833    
## I(MAFNTR^3)  5.979e-02 9.166e-02 0.655   0.516    
## G(D99)I(D93) -4.919e-04 2.050e-04 -0.999   0.316    
## G(D99)I(D93)^2 4.919e-04 2.050e-04 0.999   0.316    
## 
## Residual standard error: 160 on 84 degrees of freedom
## Multiple R-squared:  0.8472, Adjusted R-squared:  0.8399 
## F-statistic: 116.4 on 4 and 84 DF, p-value: < 2.2e-16

```

R function

$I(X^2)$ add quadratic term to the model

From the outputs, considering the scatterplot between y and X , we found that the best model to fit the data is

$$\hat{Y} = 172.4137 - 18.1073X_1 + 0.4508X_1^2$$

moreover, $R^2_{adj} = 0.805$ and RMSE= 175.6, comparing to the simple linear model, we can conclude that the quadratic model fits the data better than the simple linear regression model.

Note!

Model interpretations are not meaningful outside the range of the independent variable. Although the model appears to support the data. To make a prediction for Y , value of X should be inside the range of independent variable. Otherwise the prediction will not be meaningful.

Inclass Practice Problem 9

Suppose you want to model the monthly, y , of a product as a function of the pressure pounds per square inch (psi), at which it is produced. Four inspectors independently assign a quality score between 0 and 100 to each product, and then the quality is calculated by averaging the four scores. An experiment is conducted by varying temperature in F. Fit a second-order model to the data and sketch the scatterplot. The data are provided in PRODQUAL.csv file

Exercise 2

The amount of water used by the production facilities of a plant varies. Observations on water usage and other possibility related variables were collected for 250 months. The data are given in water.csv file. The explanatory variables are

TEMP= average monthly temperature(degree celsius)

PROD=amount of production(in hundreds of cubic)

DAYS=number of operating day in the month (days)

HOUR=number of hours shut down for maintenance (hours)

The response variable is USAGE=monthly water usage (gallons/month)

20

don't use this model

Not significant.

We discovered above that the model which model was our best?

$$\hat{y} = 8802 - 190.2[MAFNTR] + 4.587[MAFNTR]^2 - 0.0299[G(D99)I(D93)]$$

Caution: Extrapolation.

* when we predict
Stay within range
of data we collected

21