

## 8. Additional topics in linear modeling

### Outline

- We now have practical skills to
  - ① Write down linear models,
  - ② Fit them in R,
  - ③ Interpret the output in terms of parameter estimates, confidence intervals and hypothesis tests,
  - ④ Check that R is fitting the model that we intend,
  - ⑤ Check that the model we intend is appropriate for the data.
- These skills provide a foundation for many extensions helpful for particular situations.

# Topics

- The linear model formula notation in R, as a third model representation to join the subscript format and matrix format.
- Interactions between explanatory variables.
- The  $R^2$  statistic to assess model fit.
- Fitting polynomial relationships using linear models.
- Multicollinearity: What happens when two or more explanatory variables are highly correlated. How to notice it, and what to do about it.
- Power: What is the probability of rejecting the null hypothesis when the alternative is true?

# The R model formula notation

- A **formula** in `lm()` is something that looks like  $y \sim x$ .
- The R formula notation has various conventions that are designed to make it easy to specify useful models.
- `?formula` tells you everything you need to know, and more.
- The R formula for `lm()` is a way of constructing a design matrix.
- Inspect the resulting design matrix using `model.matrix()` and check you understand what R has produced. If you can do this, you can safely use the power of the formula notation.

**Question 8.1.** In a report, the model should be written in mathematical notation, not as an R formula. Why?

# Experimenting with the R formula notation

- Consider the freshman GPA data

```
gpa <- read.table("gpa.txt",header=T); head(gpa,3)
```

```
##   ID  GPA High_School ACT Year
## 1  1 0.98          61  20 1996
## 2  2 1.13          84  20 1996
## 3  3 1.25          74  19 1996
```

- We can play the game of trying out various things in R formula notation, inspecting the resulting design matrix, and figuring out how to write the model efficiently in mathematical notation.
- You can also think about whether the different models give any new insights into the data.

```
lm1 <- lm(GPA~ACT+High_School*Year,data=gpa)
coef(summary(lm1))[,1:2]
```

##	Estimate	Std. Error
## (Intercept)	-4.722613e+01	1.350854e+02
## ACT	3.708961e-02	5.946966e-03
## High_School	3.460100e-01	1.702035e+00
## Year	2.428369e-02	6.760800e-02
## High_School:Year	-1.681424e-04	8.518297e-04

- The \* here denotes inclusion of an **interaction** between High\_School and Year, written in the R output as High\_School:Year.

**Question 8.2.** Conceptually, what do you think an interaction between two variables is, and why might it be needed?

- To find out exactly what R thinks an interaction is, we can check the design matrix.

```
head(model.matrix(lm1))
```

```
##      (Intercept) ACT High_School Year High_School:Year
## 1              1  20              61 1996              121756
## 2              1  20              84 1996              167664
## 3              1  19              74 1996              147704
## 4              1  23              95 1996              189620
## 5              1  28              77 1996              153692
## 6              1  23              47 1996               93812
```

**Question 8.3.** Write out the sample model that R has computed in `lm1` using subscript notation.

# Interactions and additivity

```
lm2 <- lm(GPA~ACT+High_School+Year+High_School:Year,data=gpa)
head(model.matrix(lm2),4)
```

##	(Intercept)	ACT	High_School	Year	High_School:Year
## 1	1	20	61	1996	121756
## 2	1	20	84	1996	167664
## 3	1	19	74	1996	147704
## 4	1	23	95	1996	189620

- `lm2` has the same design matrix as `lm1`.
- We see that, in R formula notation,  $y \sim u * v$  is the same as  $y \sim u + v + u : v$ .
- In the model  $y \sim u + v$  the effects of the variables are said to be **additive**.
- In a causal interpretation of an additive model, the result of changing  $u$  to  $u_2$  and  $v$  to  $v_2$  is the sum of the marginal effect of changing  $u$  to  $u_2$  plus the marginal effect of changing  $v$  to  $v_2$ .
- The interaction term  $u : v$  breaks additivity. In this case, we can't know the consequence for the fitted value of changing  $u$  to  $u_2$  unless we know the value of  $v$ .

# The interaction between ACT and high school percentile

- We have not (yet) found any interesting effect of year. Let's drop year out of the model and look for whether there is an interaction between ACT and high school percentile for predicting freshman GPA.

```
lm3 <- lm(GPA~ACT*High_School,data=gpa)
```

**Question 8.4.** Write out the fitted sample linear model in subscript form, letting  $y_i$ ,  $a_i$ ,  $h_i$  and  $e_i$  be the freshman GPA, ACT score, high school percentile and residual error respectively for the  $i$ th student.



# Interpreting a discovered interaction

```
coef(summary(lm3))[,1:2]
```

##	Estimate	Std. Error
## (Intercept)	3.157679842	0.4788067771
## ACT	-0.046067744	0.0213355076
## High_School	-0.014405030	0.0061479608
## ACT:High_School	0.001071326	0.0002638611

**Question 8.5.** Explain in words to the admissions director what you have found about the interaction under investigation here.

## Marginal effects when there is an interaction

- Notice in 'lm3' that the coefficients for ACT score and high school percentile are negative. That is surprising!

```
ACT_centered <- gpa$ACT - mean(gpa$ACT)
HS_centered <- gpa$Hi - mean(gpa$Hi)
lm3b <- lm(GPA ~ ACT_centered * HS_centered, data = gpa)
signif(coef(summary(lm3b))[, c(1, 2, 4)], 3)
```

##	Estimate	Std. Error	Pr(> t )
## (Intercept)	2.94000	0.022900	0.00e+00
## ACT_centered	0.03640	0.005880	1.04e-09
## HS_centered	0.01190	0.001350	8.23e-18
## ACT_centered:HS_centered	0.00107	0.000264	5.46e-05

**Question 8.6.** After centering the variables, the interaction effect stays the same, but the marginal effects change sign. What is happening? Why?

## Quantifying the improvement in the model

```
s3 <- summary(lm3)$sigma
lm4 <- lm(GPA~ACT+High_School,data=gpa)
s4 <- summary(lm4)$sigma
lm5 <- lm(GPA~1,data=gpa)
s5 <- summary(lm5)$sigma
cat("s3 =",s3,"; s4 =",s4,"; s5 =",s5)

## s3 = 0.5610067 ; s4 = 0.5671605 ; s5 = 0.6345278
```

**Question 8.7.** Comment on both **statistical significance** and **practical significance** of the interaction between a prediction of freshman GPA.

# An interaction involving a factor

- Let's go back to the football field goal data.

```
goals <- read.table("FieldGoals2003to2006.csv",header=T,sep=",")
goals[1,c("Name","Teamt","FGt","FGtM1")]
```

```
##           Name Teamt  FGt FGtM1
## 1 Adam Vinatieri    NE 73.5    90
```

```
lm6 <- lm(FGt~FGtM1*Name,data=goals)
```

**Question 8.8.** What model do you think is being fitted here? Write it in subscript form, where  $y_{ij}$  is the field goal average for the  $j$ th year of kicker  $i$ , with  $i = 1, \dots, 19$  and  $j = 1, 2, 3, 4$ . Let  $e_{ij}$  be the residual error, and let  $x_{ij}$  be the previous year's average. Check your answer against the design matrix shown on the next slide.

```
X<-model.matrix(lm6) ; colnames(X)<-1:38 ; X[1:17,c(1:8,21:26)]
```

```
##      1      2 3 4 5 6 7 8      21      22      23 24 25 26
## 1    1  90.0 0 0 0 0 0 0 0.0  0.0  0.0  0  0  0
## 2    1  73.5 0 0 0 0 0 0 0.0  0.0  0.0  0  0  0
## 3    1  93.9 0 0 0 0 0 0 0.0  0.0  0.0  0  0  0
## 4    1  80.0 0 0 0 0 0 0 0.0  0.0  0.0  0  0  0
## 5    1  88.2 1 0 0 0 0 0 88.2  0.0  0.0  0  0  0
## 6    1  82.7 1 0 0 0 0 0 82.7  0.0  0.0  0  0  0
## 7    1  84.3 1 0 0 0 0 0 84.3  0.0  0.0  0  0  0
## 8    1  72.7 1 0 0 0 0 0 72.7  0.0  0.0  0  0  0
## 9    1  72.2 0 1 0 0 0 0  0.0 72.2  0.0  0  0  0
## 10   1  87.0 0 1 0 0 0 0  0.0 87.0  0.0  0  0  0
## 11   1  85.2 0 1 0 0 0 0  0.0 85.2  0.0  0  0  0
## 12   1  75.0 0 1 0 0 0 0  0.0 75.0  0.0  0  0  0
## 13   1  82.1 0 0 1 0 0 0  0.0  0.0 82.1  0  0  0
## 14   1  95.6 0 0 1 0 0 0  0.0  0.0 95.6  0  0  0
## 15   1  85.7 0 0 1 0 0 0  0.0  0.0 85.7  0  0  0
## 16   1  79.1 0 0 1 0 0 0  0.0  0.0 79.1  0  0  0
## 17   1  80.0 0 0 0 1 0 0  0.0  0.0  0.0 80  0  0
```

**Question 8.9.** Interpret the ANOVA table below.

```
anova(lm6)
```

```
## Analysis of Variance Table
```

```
##
```

```
## Response: FGt
```

```
##           Df    Sum Sq Mean Sq F value    Pr(>F)
```

```
## FGtM1       1     87.20   87.199    1.9008 0.176047
```

```
## Name       18  2252.47  125.137    2.7279 0.004565 **
```

```
## FGtM1:Name 18   417.75   23.209    0.5059 0.938592
```

```
## Residuals  38  1743.20   45.874
```

```
## ---
```

```
## Signif. codes:
```

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

# Collinear explanatory variables in a linear model

- Let  $\mathbb{X} = [x_{ij}]_{n \times p}$  be an  $n \times p$  design matrix.
- If there is a nonzero vector  $\alpha = (\alpha_1, \dots, \alpha_p)$  such that  $\mathbb{X}\alpha = \mathbf{0}$  then the columns of  $\mathbb{X}$  are **collinear**.
- Here,  $\mathbf{0}$  is the zero vector,  $(0, 0, \dots, 0)$ .
- We can write  $\mathbf{x}_j = (x_{1j}, x_{2j}, \dots, x_{nj})$  for the  $j$ th column of  $\mathbb{X}$ . Then,

$$\mathbb{X}\alpha = \alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2 + \dots + \alpha_p\mathbf{x}_p.$$

We see that  $\mathbb{X}\alpha$  can be thought of as a **linear combination of the columns of  $\mathbb{X}$** .

- Collinearity of explanatory variables has important consequences for fitting a linear model to data.
- It can also be useful to notice whether the variables are close to collinear, meaning that  $\mathbb{X}\alpha$  is small but nonzero.

## Example: an intercept with a coefficient for each factor

- Recall the mouse weight dataset. Consider a sample linear model,

$$y_{ij} = \mu + \mu_j + e_{ij}.$$

- Suppose that we don't set the  $\mu_1 = 0$  so we try to estimate both  $\mu_1$  and  $\mu_2$  at the same time as the intercept,  $\mu$ .
- Let's work with just 3 mice in each treatment group, so  $i = 1, 2, 3$  and  $j = 1, 2$ . The design matrix is therefore

```
X <- cbind(rep(1,6),rep(c(1,0),each=3),rep(c(0,1),each=3)) ; X
##      [,1] [,2] [,3]
## [1,]    1    1    0
## [2,]    1    1    0
## [3,]    1    1    0
## [4,]    1    0    1
## [5,]    1    0    1
## [6,]    1    0    1
```

- For  $\alpha = (1, -1, -1)$ , we have  $\mathbb{X}\alpha = 0$



# The least squares fit with collinear predictors

- Suppose that  $\mathbf{b}$  is a least squares coefficient vector, so that the fitted value vector  $\hat{\mathbf{y}} = \mathbb{X}\mathbf{b}$  minimizes  $\sum_{i=1}^n (y_i - \hat{y}_i)^2$ .
- Suppose that  $\mathbb{X}$  is collinear, with  $\mathbb{X}\boldsymbol{\alpha} = \mathbf{0}$ .
- Since

$$\mathbb{X}(\mathbf{b} + \boldsymbol{\alpha}) = \mathbb{X}\mathbf{b} + \mathbb{X}\boldsymbol{\alpha} = \mathbb{X}\mathbf{b} + \mathbf{0} = \mathbb{X}\mathbf{b},$$

we see that  $\mathbf{b} + \boldsymbol{\alpha}$  is also a least squares coefficient vector.

- **When  $\mathbb{X}$  is collinear, a least squares coefficient still exists, but it is not unique.**

**Question 8.10.** Let  $c$  be any number. Recall multiplication of a vector by a scalar:  $c\boldsymbol{\alpha} = (c\alpha_1, \dots, c\alpha_p)$ . Show that  $\mathbf{b} + c\boldsymbol{\alpha}$  is also a least squares fit.

## Standard errors for collinear variables

**Question 8.11.** Any variable that is part of a collinear combination of variables has infinite standard error. Why?

# What does R do if give it collinear variables?

```
mice <- read.table("femaleMiceWeights.csv",header=T,sep=",")
chow=rep(c(1,0),each=12)
hf=rep(c(0,1),each=12)
lm1 <- lm(Bodyweight~chow+hf,data=mice)
coef(summary(lm1))
```

	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	26.834167	1.039353	25.818139	6.045435e-18
## chow	-3.020833	1.469867	-2.055174	5.192480e-02

- R noticed that the three explanatory variables are collinear, and refused to fit the third

```
model.matrix(lm1)
```

##	(Intercept)	chow	hf
## 1	1	1	0
## 2	1	1	0
## 3	1	1	0
## 4	1	1	0
## 5	1	1	0
## 6	1	1	0
## 7	1	1	0
## 8	1	1	0
## 9	1	1	0
## 10	1	1	0
## 11	1	1	0
## 12	1	1	0
## 13	1	0	1
## 14	1	0	1
## 15	1	0	1
## 16	1	0	1
## 17	1	0	1
## 18	1	0	1
## 19	1	0	1



# Linearly independent vectors and matrix rank

- Columns of a matrix that are not collinear are said to be **linearly independent**.
- The **rank** of  $\mathbf{X}$  is the number of linearly independent columns.
- $\mathbf{X}$  has **full rank** if all the columns are linearly independent. In this case, we expect the least squares coefficient to be uniquely defined and so  $\mathbf{X}^T \mathbf{X}$  has non-zero determinant and is invertible.
- If  $\mathbf{X}$  does not have full rank, we can drop **linearly dependent** columns until the remaining columns are linearly independent. This is a practical approach to handling collinearity.

## Example: reducing a design matrix to full rank

```
X <- model.matrix(lm1)
```

```
det(t(X)%*%X)
```

```
## [1] 0
```

```
X2 <- X[,1:2]
```

```
det(t(X2)%*%X2)
```

```
## [1] 144
```

- Dropping the third column of  $X$  has given us a full-rank design matrix.

**Question 8.12.** The least squares fitted values are the same using the predictor matrix  $X_2$  as  $X$ . Why does dropping the last column not change the fitted values?

## Almost collinear variables

- If the determinant of  $\mathbf{X}^T\mathbf{X}$  is close to zero, the variance of the model-generated least squares coefficient vector becomes large.
- This can happen when multiple explanatory variables are included in a model which all model similar things.

**Question 8.13.** Recall our data analysis using unemployment to explain life expectancy. What would happen if we added total employment as an additional explanatory variable? (Being unemployed is not the only alternative to being employed, since only adults currently looking for work are counted as unemployed.)