

## STA302/STA1001, Weeks 10-11

Mark Ebden, 16–23 November (Section 1) and 23 November (Section 2)

With grateful acknowledgment to Alison Gibbs

## Overview

### Multiple-regression ANOVA:

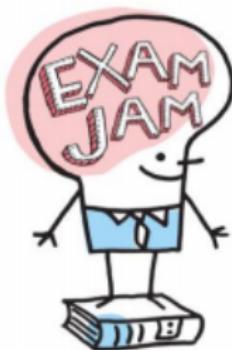
- ▶ The  $F$ -test
- ▶  $R^2$  and Adjusted  $R^2$
- ▶ Interaction terms
- ▶ ANCOVA and the partial  $F$ -test

### MLR Diagnostics



## Exam Jam

The STA302 review session will occur in SS 2135 from 10-11:30 am on 8 December. Please submit your requests for review topics closer to the time: there's a Piazza thread for this, under the 'Exam' topic.



In addition to our session: from 11 am to 3 pm there will be crafts, therapy dogs, a Photobooth, and other activities in the Sid Smith lobby. There will also be free coffee, juice, fruit, and granola bars there.

[http://www.artsci.utoronto.ca/current/exam\\_jam](http://www.artsci.utoronto.ca/current/exam_jam)

## Recap of Regression ANOVA (Week 3)

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{SST}} = \underbrace{\sum_{i=1}^n b_1^2(x_i - \bar{x})^2}_{\text{SSReg}} + \underbrace{\sum_{i=1}^n \hat{e}_i^2}_{\text{RSS}}$$

Source	SS	d.f.	MS = SS/df
Regression line	$b_1^2 S_{xx} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	1	$b_1^2 S_{xx}$
Error	$\sum_{i=1}^n \hat{e}_i^2$	$n - 2$	$S^2$
<b>Total</b>	$\sum_{i=1}^n (y_i - \bar{y})^2$	$n - 1$	-

The coefficient of determination is  $R^2 = \frac{\text{SSReg}}{\text{SST}} = 1 - \frac{\text{RSS}}{\text{SST}}, \quad 0 \leq R^2 \leq 1.$

In Weeks 9–10 we showed that the ANOVA identity can be rewritten as:

$$\underbrace{\mathbf{Y}' (\mathbf{I} - \frac{1}{n} \mathbf{J}) \mathbf{Y}}_{\text{SST}} = \underbrace{\mathbf{Y}' \left( \mathbf{H} - \frac{1}{n} \mathbf{J} \right) \mathbf{Y}}_{\text{SSReg}} + \underbrace{\mathbf{Y}' (\mathbf{I} - \mathbf{H}) \mathbf{Y}}_{\text{RSS}}$$

## Introducing Multiple-Regression ANOVA

In multiple regression, the ANOVA identity is the same as before, albeit with a different  $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ :

$$SST = SSReg + RSS$$

$$\underbrace{\mathbf{Y}' \left( \mathbf{I} - \frac{1}{n} \mathbf{J} \right) \mathbf{Y}}_{SST} = \underbrace{\mathbf{Y}' \left( \mathbf{H} - \frac{1}{n} \mathbf{J} \right) \mathbf{Y}}_{SSReg} + \underbrace{\mathbf{Y}' (\mathbf{I} - \mathbf{H}) \mathbf{Y}}_{RSS}$$

The MLR ANOVA table is similar to before, but the degrees of freedom have changed:

Source	SS	d.f.	MS = SS/df
Regression line	SSReg	$p$	$SSReg/p$
Error	RSS	$n - p - 1$	$S^2$
<b>Total</b>	<b>SST</b>	$n - 1$	—

## The $F$ -test in an MLR ANOVA table

The test hypotheses are:

- ▶  $H_0 : \beta_1 = \beta_2 = \cdots = \beta_p = 0$
- ▶  $H_a : \text{At least one of the } \beta_j \text{'s isn't 0}$

The test statistic is:

$$F_{\text{obs}} = \frac{\text{MSReg}}{\text{MSE}}$$

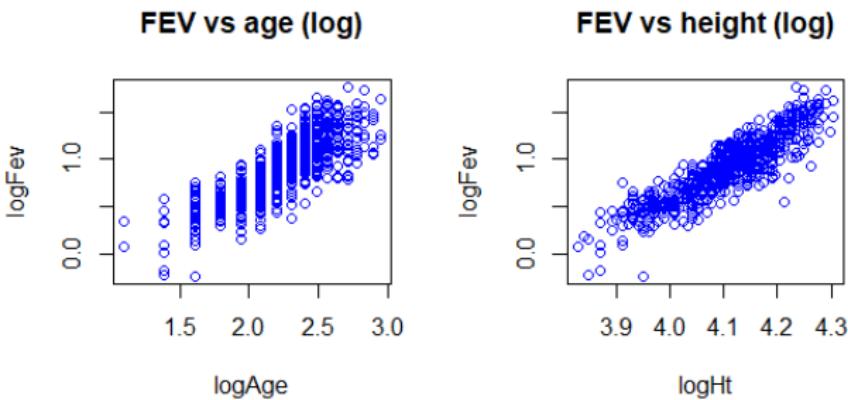
If  $H_0$  is true,  $F_{\text{obs}}$  is an observation from an  $F$  distribution with  $(p, n - p - 1)$  degrees of freedom.

- ▶ Numerator d.f.: the # of  $\beta$ 's being tested
- ▶ Denominator d.f.: the d.f. for the error

So in MLR ANOVA, we use the  $F$ -test to check for linear association between  $Y$  and any of the  $p$  predictors. If the  $F$ -test is significant, then we might ask, for which predictor(s) is there evidence of a linear association with  $Y$ ? Some pitfalls in answering this question are investigated in Chapter 7.

## Example of an *F*-test: the fev database

```
a2 = read.table("DataA2.txt",sep=" ",header=T) # Load the data set
logFev <- log(a2$fev); logAge <- log(a2$age); logHt <- log(a2$ht)
par(mfrow=c(1,2))
plot(logAge,logFev,type="p",col="blue",pch=21, main="FEV vs age (log)")
plot(logHt,logFev,type="p",col="blue",pch=21, main="FEV vs ht (log)")
mod1 = lm(logFev~logAge+logHt)
```



## SLR in the fev database

```
##  
## Call:  
## lm(formula = logFev ~ logAge)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -0.60857 -0.13532  0.00227  0.14329  0.56348  
##  
## Coefficients:  
##                 Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -0.98772    0.05756 -17.16   <2e-16 ***  
## logAge       0.84615    0.02535  33.38   <2e-16 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.2026 on 652 degrees of freedom  
## Multiple R-squared:  0.6309, Adjusted R-squared:  0.6303  
## F-statistic: 1114 on 1 and 652 DF,  p-value: < 2.2e-16
```

## SLR in the fev database

```
##  
## Call:  
## lm(formula = logFev ~ logHt)  
##  
## Residuals:  
##      Min       1Q   Median      3Q     Max  
## -0.69369 -0.09122  0.01145  0.09832  0.44965  
##  
## Coefficients:  
##                 Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -11.92110    0.25577 -46.61 <2e-16 ***  
## logHt        3.12418    0.06223  50.20 <2e-16 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.1512 on 652 degrees of freedom  
## Multiple R-squared:  0.7945, Adjusted R-squared:  0.7941  
## F-statistic: 2520 on 1 and 652 DF,  p-value: < 2.2e-16
```

## MLR in the fev database

```
##  
## Call:  
## lm(formula = logFev ~ logAge + logHt)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -0.62020 -0.08894  0.01166  0.09807  0.46645  
##  
## Coefficients:  
##                 Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -10.29520    0.39196 -26.266 < 2e-16 ***  
## logAge       0.18045    0.03346   5.392 9.74e-08 ***  
## logHt        2.62968    0.11010  23.884 < 2e-16 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.1481 on 651 degrees of freedom  
## Multiple R-squared:  0.8033, Adjusted R-squared:  0.8026  
## F-statistic: 1329 on 2 and 651 DF,  p-value: < 2.2e-16
```

## $R^2$ for MLR ANOVA

Let's consider the coefficient of determination for MLR ANOVA, a.k.a. the "coefficient of **multiple** determination":

$$R^2 = \frac{SSReg}{SST} = \frac{\mathbf{Y}' (\mathbf{H} - \frac{1}{n} \mathbf{J}) \mathbf{Y}}{\mathbf{Y}' (\mathbf{I} - \frac{1}{n} \mathbf{J}) \mathbf{Y}}$$

It's not the square of correlation  $r$  anymore! Correlation is between two variables, whereas we have potentially many variables now.

However, as before, it's the proportion of the total sample variability in the  $Y$ 's explained by the regression model.

**Question:** What happens to  $R^2$  when you add more predictor variables?

## The effect on $R^2$ of additional predictors

Each time a predictor variable is added, SST stays the same because it depends on  $\mathbf{Y}$  only.

However, adding a new predictor variable often improves (decreases) RSS: a richer model will often lead to a better fit, i.e. less error. Recall that  $\text{RSS} = \hat{\mathbf{e}}' \hat{\mathbf{e}}$ . A least-squares minimization of RSS, with additional predictors now, is minimizing over a larger-dimensional space. This guarantees that the minimum is at least as small. So, at worst, RSS will stay the same (if you add a predictor that's ignored by fitting  $\hat{\beta}_j = 0$ ), and usually it will get better.

If  $\text{SST}$  is constant and  $\text{RSS}$  decreases,  $\text{SSReg}$  must increase. Therefore  $R^2$  will increase. (Put another way, the  $\mathbf{H}$  in the numerator will have changed.)

## Adjusted $R^2$

Because  $R^2$  generally increases with the number of predictors, how do we compare the  $R^2$  for a simple model to the  $R^2$  for a many-variable model?

We can use the **Adjusted  $R^2$** , a better measure of the model fit. It is adjusted for the number of predictors in the model.

$$\text{Adj } R^2 = 1 - (n - 1) \frac{\text{MSE}}{\text{SST}} = 1 - \frac{n - 1}{n - p - 1} \frac{\text{RSS}}{\text{SST}}$$

With additional predictor variables, the Adjusted  $R^2$  will only increase if MSE decreases.



## Adjusted $R^2$ in action: First, reviewing regression ANOVA

For the fev vs age SLR dataset (HW2, question 1),  $n = 654$  and  $p = 1$ .

From Weeks 9–10 slide 18,  $R^2 \approx 0.5722$  and Adj  $R^2 \approx 0.5716 \approx R^2$ , a difference of approximately only 0.1%.

Taking logs, and rerunning the analysis, today we got  $R^2 \approx 0.6309$  and Adj  $R^2 \approx 0.6303 \approx R^2$ .

## Adjusted $R^2$ in action: MLR ANOVA

Let's compare the (adjusted) coefficients of determination for a small dataset, with and without an extra predictor.

Consider just the first ten points in the fev database (A = abridged):

```
set.seed(1)
N<-10; u <- sample(length(logFev),N)
logFevA<-logFev[u]; logAgeA<-logAge[u]
rA<-rnorm(N) # A new potential predictor

mod2 = lm(logFevA~logAgeA)
mod3 = lm(logFevA~logAgeA+rA)
summary(mod2) # SLR ANOVA
summary(mod3) # MLR ANOVA
```

Note that rA is noise, but adding it still increases the  $R^2$ .

## Results of SLR ANOVA

```
##  
## Call:  
## lm(formula = logFevA ~ logAgeA)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -0.34977 -0.04767 -0.00790  0.10280  0.26091  
##  
## Coefficients:  
##                 Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -1.6288     0.5944 -2.740  0.02544 *  
## logAgeA      1.1232     0.2523  4.452  0.00213 **  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.1747 on 8 degrees of freedom  
## Multiple R-squared:  0.7125, Adjusted R-squared:  0.6765  
## F-statistic: 19.82 on 1 and 8 DF,  p-value: 0.002132
```

## Results of MLR ANOVA

```
##  
## Call:  
## lm(formula = logFevA ~ logAgeA + rA)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -0.32561 -0.05576 -0.01012  0.05902  0.29785  
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -1.72678    0.64144 -2.692  0.03099 *  
## logAgeA     1.16367    0.27176  4.282  0.00365 **  
## rA          0.03408    0.05727  0.595  0.57055  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.1822 on 7 degrees of freedom  
## Multiple R-squared:  0.7263, Adjusted R-squared:  0.6481  
## F-statistic: 9.289 on 2 and 7 DF,  p-value: 0.01072
```

## Overview

Multiple-regression ANOVA:

- ▶ The  $F$ -test
- ▶  $R^2$  and Adjusted  $R^2$
- ▶ **Interaction terms**
- ▶ ANCOVA and the partial  $F$ -test

MLR Diagnostics



## Regression model with interaction

An *additive* model (**no interaction**):

$$\text{fev} = \beta_0 + \beta_1 \text{age} + \beta_2 \text{ht} + e$$

A model that is *not* additive (**has an interaction term**):

$$\text{fev} = \beta_0 + \beta_1 \text{age} + \beta_2 \text{ht} + \beta_3 \text{age} \times \text{ht} + e$$

It can help us answer the question, “Does the relationship of `fev` with `age` depend on `height`? ”

Two explanatory variables are said to *interact* if the effect that one of them has on the response depends on the value of the other.

How can we quantitatively assess this?

## MLR ANOVA without interaction

```
##  
## Call:  
## lm(formula = logFev ~ logAge + logHt)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -0.62020 -0.08894  0.01166  0.09807  0.46645  
##  
## Coefficients:  
##                 Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -10.29520    0.39196 -26.266 < 2e-16 ***  
## logAge       0.18045    0.03346   5.392 9.74e-08 ***  
## logHt        2.62968    0.11010  23.884 < 2e-16 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.1481 on 651 degrees of freedom  
## Multiple R-squared:  0.8033, Adjusted R-squared:  0.8026  
## F-statistic: 1329 on 2 and 651 DF,  p-value: < 2.2e-16
```

## MLR ANOVA with interaction

```
##  
## Call:  
## lm(formula = logFev ~ logAge * logHt)  
##  
## Residuals:  
##       Min        1Q     Median        3Q       Max  
## -0.64913 -0.08337  0.01099  0.09729  0.42260  
##  
## Coefficients:  
##             Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -4.5057    1.5322 -2.941 0.003392 **  
## logAge      -2.4648    0.6781 -3.635 0.000300 ***  
## logHt        1.2039    0.3809  3.160 0.001649 **  
## logAge:logHt  0.6495    0.1663  3.906 0.000104 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.1465 on 650 degrees of freedom  
## Multiple R-squared:  0.8078, Adjusted R-squared:  0.8069  
## F-statistic: 910.4 on 3 and 650 DF,  p-value: < 2.2e-16
```

## Considering the *t*-test result

We called `lm(logFev~logAge*logHt)`, which is equivalent to calling  
`lm(logFev~logAge+logHt+logAge:logHt)`

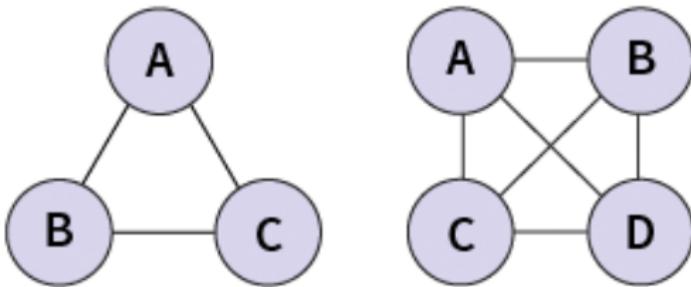
From the *t*-test regarding `logAge:logHt`, we can conclude that we have evidence that the coefficient of  $\text{age} \times \text{ht}$  is statistically significantly different from 0, given that the other terms are in the model.

Note that this model has a slightly smaller MSE and larger Adj  $R^2$  than the additive model.

We can conclude that adding the interaction term is worthwhile.

Should we routinely add interaction terms? (Hint: consider combinatorics.)

## When to add interaction terms



When to add them can also be considered a research question.

However, a standard practice is that if an interaction term is in the model, we also include the individual terms for the predictor variables, even if their coefficients are not statistically significantly different from 0.

## Next steps

- ▶ Try Chapter 5's **question 2**
- ▶ Remember that on Tuesday **21 November** we'll start at 11:10 am
- ▶ Solutions to Chapter 5's question 1 will be uploaded by 23 November



## Appendix



- ▶ ANCOVA
- ▶ The partial  $F$ -test
- ▶ MLR Diagnostics

## What happens when we add a $(\log \text{Height})^2$ term and all interactions?

```
##  
## Call:  
## lm(formula = logFev ~ logAge * logHt * logH2)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -0.66838 -0.08213  0.00931  0.09914  0.41712  
##  
## Coefficients:  
##                               Estimate Std. Error t value Pr(>|t|)  
## (Intercept)           -1960.63    1760.15 -1.114   0.266  
## logAge                  979.32     744.44  1.316   0.189  
## logHt                 1425.69    1310.11  1.088   0.277  
## logH2                  -345.49     324.95 -1.063   0.288  
## logAge:logHt            -714.54     550.66 -1.298   0.195  
## logAge:logH2              173.53     135.77  1.278   0.202  
## logHt:logH2                27.91     26.86  1.039   0.299  
## logAge:logHt:logH2      -14.02      11.16 -1.257   0.209  
##  
## Residual standard error: 0.1465 on 646 degrees of freedom  
## Multiple R-squared:  0.8088, Adjusted R-squared:  0.8067  
## F-statistic: 390.4 on 7 and 646 DF,  p-value: < 2.2e-16
```

## Analysis of the multi-parameter model

The  $F$ -test had a null hypothesis that  $\beta_1 = \dots = \beta_7 = 0$ , which was rejected with a  $p$ -value of approximately  $10^{-16}$ .

However, the  $p$ -values for the seven tests were all well above 0.05. Namely,  $H_0 : \beta_j = 0$  was never rejected. Do we conclude that all of the  $\beta_j$ 's are probably zero?

Answer: No, each  $t$ -test is for the effect of one explanatory variable given that the others are in the model.

## Second example: A meadowfoam experiment



Meadowfoam is a flower found on the West Coast, which produces an oil of use to the cosmetics and hair-care industries.

A randomized experiment was conducted to explore the effect of growing conditions on the number of flower blooms per plant.

## Example 2: dataset overview

There were  $6 \times 2 = 12$  unique treatments, for:

- ▶ Six light intensities: Intensity  $\in \{150, 300, 450, 600, 750, 900\}$ , measured in  $\mu\text{mol}/\text{m}^2/\text{s}$
- ▶ Two timings at which light began: Time is 1 if late, 2 if early

Each treatment was applied in two trials, so there were 24 trials in total.

The response variable,  $Y$ , known as Flowers in the dataset, was the average number of flowers observed per plant (across ten plants in a single pot).

Two questions of interest: What's the effect on the number of flowers per plant, of:

- ▶ Timing
- ▶ Light intensity

A scientific paper with background, as optional reading:

<http://agris.fao.org/agris-search/search.do?recordID=US9500398>

## The data

```
library(Sleuth3)
print(case0901)
```

	Flowers	Time	Intensity
## 1	62.3	1	150
## 2	77.4	1	150
## 3	55.3	1	300
## 4	54.2	1	300
## 5	49.6	1	450
## 6	61.9	1	450
## 7	39.4	1	600
## 8	45.7	1	600
## 9	31.3	1	750
## 10	44.9	1	750
## 11	36.8	1	900
## 12	41.9	1	900
## 13	77.8	2	150
## 14	75.6	2	150
## 15	69.1	2	300
## 16	78.0	2	300
## 17	57.0	2	450
## 18	71.1	2	450
## 19	62.9	2	600

## Strategy

We'll set categorical variable  $t$  to be 0 or 1 for late and early, respectively.

We'll also treat Intensity as a *categorical* variable (!)

We can do this because Intensity has a small number of values with multiple observations for each. This approach may be useful for learning which intensity leads to the highest value of response variable, without imposing a particular form of relationship on Intensity versus Flowers. It may be linear, quadratic, etc. (You may wish to examine slides 71–72.)

Shall we define six new indicator variables?

$$i_{150} = \begin{cases} 1 & \text{if Intensity} = 150 \\ 0 & \text{otherwise} \end{cases} \quad \dots \quad i_{900} = \begin{cases} 1 & \text{if Intensity} = 900 \\ 0 & \text{otherwise} \end{cases}$$

## Economical representation

Using all six indicator variables is redundant: e.g. if five variables are zero, you know that the sixth is 1. In the  $24 \times 8$  design matrix, the columns for  $i150, \dots, i900$  contain a linear dependence.

This will lead to an error in R when running the `lm` command.

In general: For a categorical variable with  $k$  categories, you need  $k - 1$  indicator variables.



The model will be:

$$Y = \beta_0 + \beta_1 i150 + \beta_2 i300 + \beta_3 i450 + \beta_4 i600 + \beta_5 i750 + \beta_6 t + e$$

## Results

This code ensures that Intensity = 900 will be the reference level, as it's listed first:

```
i <- factor(case0901$Intensity, levels=c(900,150,300,450,600,750))
myFit <- lm(Flowers ~ i + as.factor(Time), data=case0901)
summary(myFit)
```

The fitted model is:

$$\hat{Y} = 37.8 + 29.4i150 + 20.2i300 + 16.0i450 + 6.1i600 + 1.6i750 + 12.2t$$

When Intensity is 150, and Time is early, what's the estimate of the mean number of flowers per plant?

What does the intercept estimate?

```

## 
## Call:
## lm(formula = Flowers ~ i + as.factor(Time), data = case0901)
## 
## Residuals:
##    Min     1Q Median     3Q    Max 
## -8.979 -4.308 -1.342  5.204 10.204 
## 
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 37.846    3.629   10.430 8.33e-09 ***
## i150        29.350    4.751    6.178 1.01e-05 ***
## i300        20.225    4.751    4.257 0.000532 ***
## i450        15.975    4.751    3.362 0.003697 **  
## i600        6.125     4.751    1.289 0.214601    
## i750        1.600     4.751    0.337 0.740415    
## as.factor(Time)2 12.158    2.743    4.432 0.000365 ***
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## 
## Residual standard error: 6.719 on 17 degrees of freedom
## Multiple R-squared:  0.8231, Adjusted R-squared:  0.7606 
## F-statistic: 13.18 on 6 and 17 DF,  p-value: 1.427e-05

```

## Is timing important?



Let's consider  $H_0 : \beta_6 = 0$  versus  $H_a : \beta_6 \neq 0$ . The test statistic is about 4.43, with a  $p$ -value calculated from a  $t_{17}$  distribution of about 0.0004.

**Yes**, timing is important. After accounting for the effect of intensity, there is strong evidence that the mean of the number of flowers per plant differs between the early and late timings.

Holding intensity constant, we get on average 12.2 flowers per plant more with early timing.

## Is intensity important?



Let's consider  $H_0 : \beta_1 = 0$  versus  $H_a : \beta_1 \neq 0$ . The  $p$ -value is less than 0.0001, so there is strong evidence that  $\beta_1 \neq 0$  given other variables in the model.

If  $\beta_1 = 0$ , the model would be the same for intensities of 150 and 900.

So **yes** intensity is important. We conclude that a light intensity of 150 gives, on average, a different number of flowers per plant than an intensity of 900.

*Individual tests for  $\beta_1, \dots, \beta_5$  compare the mean response at a certain intensity to that at an intensity of 900.*

## Is intensity important?

What we really want to test is  $H_0 : \beta_1 = \dots = \beta_5 = 0$  versus  $H_a$  : at least one of  $\beta_1, \dots, \beta_5$  isn't zero.

We should run a **partial F-test**. This tests whether a subset of  $\beta$ 's are zero simultaneously.

The approach is:

1. Fit the model with all predictor variables (known as the *full model*), and calculate RSS, known as  $\text{RSS}(\text{full})$
2. Fit the model without the predictor variables whose coefficients we're testing (known as the *reduced model*), and calculate RSS, known as  $\text{RSS}(\text{reduced})$
3. Calculate the observed  $F$ :

$$F = \frac{(\text{RSS}(\text{reduced}) - \text{RSS}(\text{full})) / (\text{df}_{\text{reduced}} - \text{df}_{\text{full}})}{\text{RSS}(\text{full}) / \text{df}_{\text{full}}}$$

## The partial $F$ -test

We know that:

- ▶ RSS in reduced model  $\geq$  RSS in full model
- ▶ SSReg in reduced model  $\leq$  SSReg in full model
- ▶ SST in reduced model = SST in full model

Note that  $df_{full}$  is the number of degrees of freedom in the error for the full model. The difference  $df_{reduced} - df_{full}$  is the number of parameters that you're testing in the partial  $F$ -test.

It can be shown that, under  $H_0$ ,  $F_{obs}$  has an  $F$  distribution with  $(df_{reduced} - df_{full}, df_{full})$  degrees of freedom.

The **intuition** behind the test is: Did RSS go down by a statistically significant amount when new predictors were added to the model?

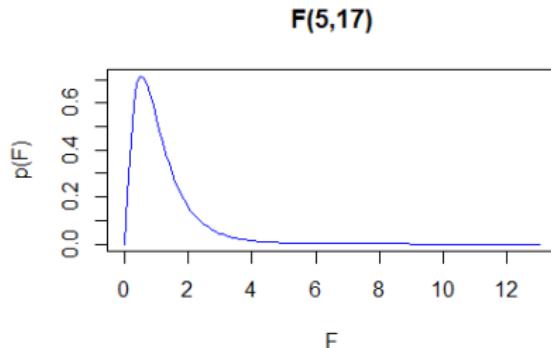
- ▶ Equivalently: Did  $R^2$  increase by a statistically significant amount?

## Back to our example: Is intensity important?

$H_0 : \beta_1 = \dots = \beta_5 = 0$  versus  $H_a : \text{at least one of } \beta_1, \dots, \beta_5 \text{ isn't zero.}$

We obtain a test statistic of

$$F_{\text{obs}} \approx \frac{(3451 - 767)/5}{767/17} \approx 11.9$$



There is strong evidence that not all of  $\beta_1, \dots, \beta_5$  are zero, given that time is in the model. So we have reconfirmed that **yes** intensity is important.

## The ANOVA table for the Meadowfoam dataset

We have decomposed  $SS_{\text{Reg}}$  into two components: intensity and timing.

Source	df	SS	MS	F
Regr(timing)	1	887	887	$887/45.15 = 19.6$
Regr(intensity)	5	2684	538	$538/45.15 = 11.9$
Error	17	767	45.15	
Total	23	4338		

Note that  $887/45.15 \approx 19.6 \approx (4.43)^2$ .

Also note that we could carry out a partial  $F$ -test on  $H_0 : \beta_j = 0$  versus  $H_a : \beta_j \neq 0$ , i.e. on one parameter. Of course, this assumes all other variables are in the model.

**Exercise:** Try this for  $\beta_5$ , i.e. for  $i750$ 's coefficient.

## New question

Does the way light intensity affects the mean of the number of flowers per plant depend on timing?



In setting up a model to answer this question, we'll continue to model timing as a qualitative variable, but we'll begin to model intensity as a quantitative variable. (You may wish to examine slides 71–72.)

## Analysis of Covariance (ANCOVA)

In ANCOVA, the predictors include both quantitative variables and qualitative variables, e.g.  $d \in \{0, 1\}$ .

**Parallel regression lines:**

$$Y = \beta_0 + \beta_1 x + \beta_2 d + e$$

**Regression lines with equal intercepts but different slopes:**

$$Y = \beta_0 + \beta_1 x + \beta_3 d x + e$$

**Unrelated regression lines:**

$$Y = \beta_0 + \beta_1 x + \beta_2 d + \beta_3 d x + e$$

The last cases are examples of introducing an *interaction*, as we saw earlier.

## Using ANCOVA to answer our new question

$$Y = \beta_0 + \beta_1 x + \beta_2 d + \beta_3 dx + e$$

We'll test whether the resulting change in  $Y$  (Flowers) when  $x$  (Intensity) changes is the same for early- versus late timings ( $d = 1$  or  $0$ ). In other words,  $H_0 : \beta_3 = 0$ .

This isn't the same as asking: "Is the relationship between  $Y$  and Intensity the same for early and late timings? Do they have the same line?" (What is the hypothesis test in that case?)

The R code for our test is:

```
myFit <- lm(Flowers ~ Intensity * as.factor(Time), data=case0901)
summary(myFit)
```

## R output

```
##  
## Call:  
## lm(formula = Flowers ~ Intensity * as.factor(Time), data = case0901)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -9.516 -4.276 -1.422  5.473 11.938  
##  
## Coefficients:  
##                               Estimate Std. Error t value Pr(>|t|)  
## (Intercept)                71.623333   4.343305 16.491 4.14e-13 ***  
## Intensity                 -0.041076   0.007435 -5.525 2.08e-05 ***  
## as.factor(Time)2           11.523333   6.142360  1.876  0.0753 .  
## Intensity:as.factor(Time)2  0.001210   0.010515  0.115  0.9096  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 6.598 on 20 degrees of freedom  
## Multiple R-squared:  0.7993, Adjusted R-squared:  0.7692  
## F-statistic: 26.55 on 3 and 20 DF,  p-value: 3.549e-07
```

## Conclusions regarding the interaction

From the unrelated-regressions model, there is no evidence that the effect of light intensity on the number of flowers per plant differs with timing ( $p \approx 0.91$ ).

If there were significant interactions (as we saw in the `fev` example), it would be difficult to talk about the effects of the individual predictor variables because they'd depend on the value of others.

**Next step:** Since the coefficient of interaction is not statistically significantly different from 0, remove it so that we can talk about the individual effects of timing and intensity.

## R output

```
##  
## Call:  
## lm(formula = Flowers ~ Intensity + as.factor(Time), data = case0901)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -9.652 -4.139 -1.558  5.632 12.165  
##  
## Coefficients:  
##                               Estimate Std. Error t value Pr(>|t|)  
## (Intercept)           71.305833   3.273772 21.781 6.77e-16 ***  
## Intensity            -0.040471   0.005132 -7.886 1.04e-07 ***  
## as.factor(Time)2    12.158333   2.629557  4.624 0.000146 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 6.441 on 21 degrees of freedom  
## Multiple R-squared:  0.7992, Adjusted R-squared:  0.78  
## F-statistic: 41.78 on 2 and 21 DF,  p-value: 4.786e-08
```

## Continued analysis

There is strong evidence ( $p < 0.001$ ) that light intensity affects the number of flowers per plant over and above timing.

For a given timing, increasing the light intensity by  $100 \mu\text{mol}/\text{m}^2/\text{s}$  decreases the number of flowers per plant on average by approximately 4.0.

**Exercise:** Show that the 95% CI for this decrease is  $(-5.1, -3.0)$ .

There is strong evidence ( $p \approx 0.0001$ ) that timing affects the number of flowers per plant over and above light intensity.

For a given intensity, introducing early timing increases the number of flowers per plant on average by approximately 12.2.

**Exercise:** Show that the 95% confidence interval for this increase is  $(6.7, 17.6)$ .

## Continued analysis

We could have fit two separate regression lines by splitting the data into the twelve early observations and the twelve late observations.

Advantages of using ANCOVA included:

- ▶ We have tests for equal slopes and intercepts
- ▶ We have higher  $df_{\text{error}}$ , meaning the power increases and the CIs are narrower
- ▶ We get a better estimate of the error variance based on 24 observations rather than 12

A possible disadvantage of using ANCOVA was:

- ▶ An implicit assumption that both groups have the same error variance

## Recap of SLR diagnostics: The Seven C's (7 Checks)

1. Plot *standardized residuals* to help determine whether the proposed regression model is a valid model
2. Identify any *leverage points*
3. Identify any *outliers*
4. Identify any *influential points*
5. Assess the assumption of *error homoscedasticity*
6. For time series: examine whether the data are *correlated over time*
7. Assess the assumption of *normal errors*



## Are the Seven Seas Unique?

Ya'qubi was a historian living in ninth-century Asia and Africa. He wrote, "Whoever wants to go to China must cross seven seas, each one with its own colour, wind, fish and breeze." He identified:

Persian Gulf  
Arabian Sea  
Bay of Bengal

Strait of Malacca  
Singapore Strait

Gulf of Thailand  
South China Sea



Optional material

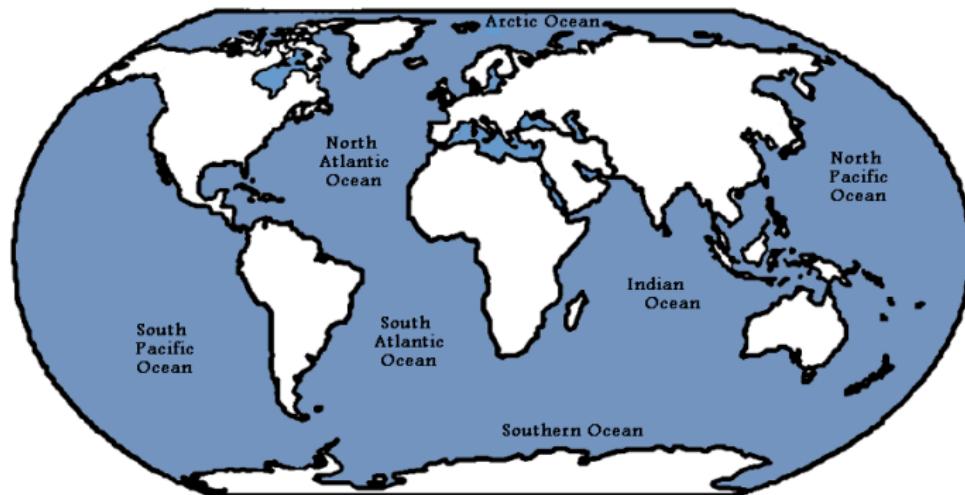
## Alternative Seven Seas

Other lists of Seven Seas date from 2300 B.C. to the present, e.g.:

Pacific Ocean  
Atlantic Ocean  
Indian Ocean

Arctic Ocean  
Mediterranean Sea

Caribbean Sea  
Gulf of Mexico



Optional material

## MLR diagnostics: The Seven C's (7 Checks) revisited

Chapter 6 of our textbook lists seven C's for MLR. Five of them are similar to those for SLR:

1. Plot *standardized residuals* to help determine whether the proposed regression model is a valid model. Do further analysis, not covered here.
2. Identify any *leverage points*
3. Identify any *outliers*
6. Assess the assumption of *error homoscedasticity*
7. For time series: examine whether the data are *correlated over time*

But, two checks are new, as they don't apply to SLR:

4. Assess the effect of each  $X$  on  $Y$
5. Assess *multicollinearity*

## Five traditional C's, now applied to MLR

1. The  $i$ th standardized residual is, as before,

$$r_i = \frac{\hat{e}_i}{S\sqrt{1 - h_{ii}}}$$

However,  $S = \sqrt{\text{RSS}/(n - p - 1)}$  is the MLR estimate of  $\sigma$ ; more important, there are superior techniques available, beyond the scope of this course.

2. To find leverage, we compute  $\mathbf{H}$  as per earlier lectures. The threshold is:

$$h_{ii} > \frac{2(p + 1)}{n}$$

3. As before, outliers have  $|r_i| > 2$  depending on the size of the dataset.
6. The constancy of variance is checked in a new way for MLR (see Check 1).
7. Correlations over time can be assessed as they were for SLR.

## What about the two missing C's?

In MLR, we should still assess the assumption of **normal errors**: for this course, we can use the normal quantile plot of residuals as with SLR.

We can also still identify any **influential points**. The influence statistics are the same as those for simple linear regression:

$$\text{DFBETA}_{ik} = \frac{\hat{\beta}_k - \hat{\beta}_{k(i)}}{\text{s.e. of } \hat{\beta}_{k(i)}} \quad \text{DFFITS}_i = \frac{\hat{y}_i - \hat{y}_{i(i)}}{\text{s.e. of } \hat{y}_{i(i)}}$$

$$D_i = \frac{\sum_{j=1} (\hat{y}_{j(i)} - \hat{y}_j)^2}{2S^2} = \frac{r_i^2 h_{ii}}{2(1 - h_{ii})}$$

However, the threshold formulae are generalized for MLR as follows:

- ▶ DFBETAS  $> 2/\sqrt{n}$
- ▶ DFFITS  $> 2\sqrt{\frac{p+1}{n}}$
- ▶ Cook's distance  $D > \frac{4}{n-(p+1)}$

In each case, some authors check for a large gap as well, just as with SLR

## Focus on MLR Check 1: Residual plots

Plotting (standardized) residuals versus  $X_j$  for  $j \in \{1, \dots, p\}$  helps us to look for:

- ▶ Curvature
- ▶ Influential points
- ▶ Outliers

Plotting (standardized) residuals versus  $Y$  helps us to look for:

- ▶ Nonconstant variance
- ▶ Outliers

Plotting residuals versus other potential predictors can help expand our model as appropriate, as mentioned in our SLR work.

And as with SLR, residual plots are not the only way to assess model assumptions. For example, plots of  $Y$  vs  $X_j$  help us answer:

- ▶ Is the linear model appropriate?
- ▶ Are there unusual points? (e.g. potential outliers or influential points)

We can also look at the added variable plots (Check 4, to come).

## MLR Check 5: Multicollinearity

**Multicollinearity** occurs when there is lots of correlation among the  $X$ 's. We use the term interchangeably with “ill-conditioning”.

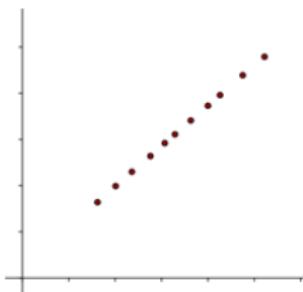
When explanatory variables are highly correlated, it's difficult or impossible to measure the individual variable's influence on the response.

The fitted equation is unstable:

- ▶ The estimated regression coefficients vary widely from data set to data set (even if the data sets are very similar), and depending on which other predictor variables are in the model
- ▶ An estimated coefficient may have opposite sign to what you'd expect
- ▶ A coefficient might not be statistically significantly different from zero even though there is a strong relationship between the  $X$  and  $Y$  when only considering  $X$  and  $Y$

## Multicollinearity

Recall:  $\hat{\beta} = \mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ . If some the  $X$ 's are perfectly correlated,  $\mathbf{X}'\mathbf{X}$  is singular and we can't calculate  $\mathbf{b}$ .



Put another way, in terms of the Week 8 SLR material: if  $X$  contains linearly dependent columns,  $X$  has a rank below  $p + 1$ . Therefore  $(\mathbf{X}'\mathbf{X})^{-1}$  has a rank below  $p + 1$ . A matrix must have full rank to be invertible.

In the case of  $\mathbf{X}'\mathbf{X}$  close to singular, the determinant of  $\mathbf{X}'\mathbf{X}$  will be near 0. Therefore,  $\text{var}(\mathbf{b}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$  will be large. This means that the standard errors of the estimated coefficients will be large, so we'll have "inefficient" estimates. (We can't make precise statements about their values.)

## Quantifying Multicollinearity

Let  $R_j^2$  represent the coefficient of multiple determination obtained when the  $j$ th predictor variable is regressed against the other predictor variables.

The **variance inflation factor** is

$$\text{VIF}_j = \frac{1}{1 - R_j^2}$$

A large  $\text{VIF}_j$  is a sign of multicollinearity. Rules of thumb:

- ▶ If  $5 \lesssim \text{VIF}_j \lesssim 10$ , the effects of multicollinearity might be seen (this is a warning)
- ▶ If  $\text{VIF}_j \gtrsim 10$ , there is a serious problem

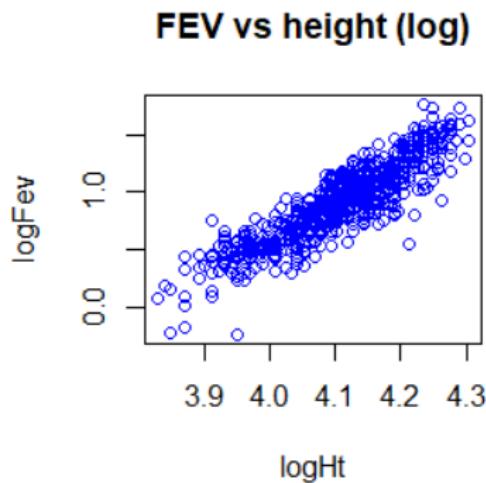
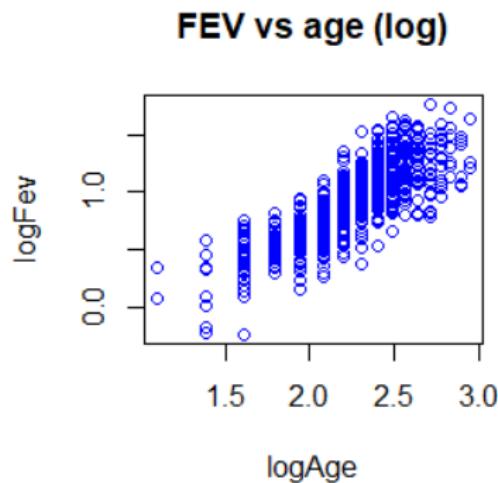
However, if the  $j$ th variable's coefficient is statistically significantly different from zero, this can be an indication not to worry as much about a high  $\text{VIF}_j$ .

Optional material: **Tolerance** is defined as  $1/\text{VIF}_j$

## Multicollinearity in the FEV dataset

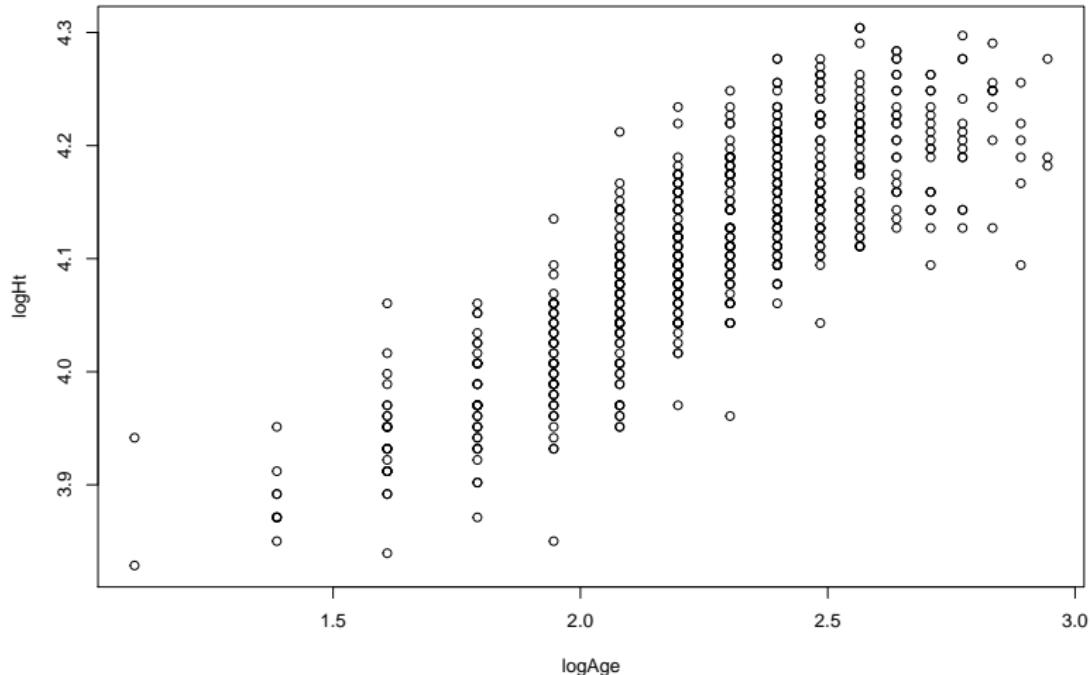
We'd calculated (on slides 10 and 20) that  $\beta_1 \approx 0.18$  for the relationship between logFEV and logAge in a particular MLR context.

Here are the graphs from slide 7:



For this dataset,  $VIF_1 = VIF_2 \approx 3.3$ .

## Scatterplot of the predictor variables



## Proposed solution 1 to multicollinearity: Ridge regression

When fitting regression models with serious multicollinearity, you may try ridge regression:  $\hat{\beta} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}'\mathbf{Y}$



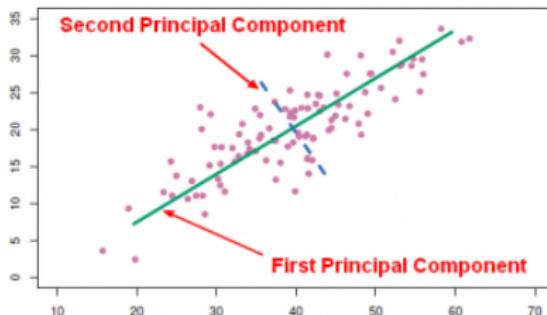
$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

- ▶ “Beefing up” the diagonal of the matrix being inverted makes the problem better-conditioned
- ▶ In Week 12, if time permits we might discuss how to choose  $\lambda$  (it would be optional material)

## Proposed solution 2: PCR

In principal component regression,  $\hat{\beta} = (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{Y}$

- $\mathbf{Z}$  is a lower-dimensional version of  $\mathbf{X}$ , e.g.  $n \times 2$  instead of  $n \times 3$
- Suppose  $\mathbf{X}$  consists of two predictor variables that are noticeably correlated, and you run PCA (principal component analysis) to discover the direction of maximum variation (along the first principal component, a  $2 \times 1$  vector):



- Projecting the corresponding two  $n$ -columns of  $\mathbf{X}$  onto the principal component results in a single  $n$ -column (distances along the teal line)
- This  $n$ -column goes into forming the matrix  $\mathbf{Z}$  (dimension  $n \times 2$  here)
- You then run SLR on  $\mathbf{Z}$  which has just one hybridized predictor variable

## MLR Check 4: Added variable plots

An **added variable plot** allows you to *visualize* the relationship between a response variable and an explanatory variable over and above the other explanatory variables.

Such plots are also known as *partial regression plots*, *adjusted variable plots*, or *partial residual plots*.

The technique is based on the concept of **partial correlation**. The partial correlation between  $X_1$  and  $X_2$  given a set of  $n$  other variables

$Z = \{Z_1, Z_2, \dots, Z_n\}$  is the correlation between:

- ▶ The residuals  $\hat{e}_{X_1}$  resulting from the linear regression of  $X_1$  versus  $Z$
- ▶ The residuals  $\hat{e}_{X_2}$  resulting from the linear regression of  $X_2$  versus  $Z$

## Added variable plots

For the  $j$ th added variable plot, we first divide  $X$  into  $X_j$  and the other  $X$ 's (excluding  $X_j$ ). Then we plot:

- ▶  $x$ -axis: Residuals from the regression of  $X_j$  versus the other  $X$ 's
- ▶  $y$ -axis: Residuals from the regression of  $Y$  versus the other  $X$ 's

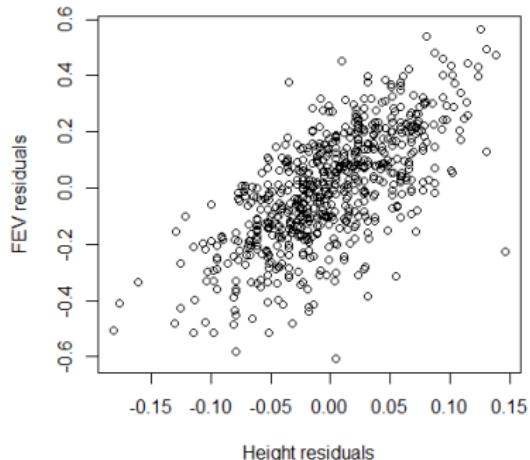
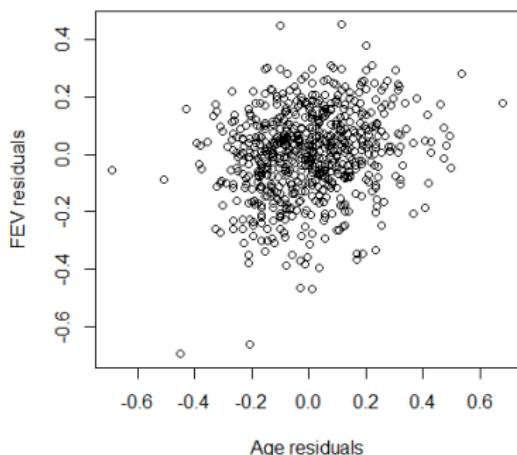
The correlation between the  $x$ - and  $y$ -axes here is a famous example of a partial correlation.

In an added variable plot, a linear pattern indicates that  $X_j$  is useful in the model over and above the other explanatory variables. In other words, the plot shows the strength of linear relationship  $Y$  and  $X_j$  over and above other variables.

The plot is also useful for detecting nonlinear relationships (polynomial in  $X_j$  for example), outliers, nonconstant variance, and influential points.

## Example: FEV data

```
par(mfrow=c(1,2))
xAxis <- lm(logAge ~ logHt); yAxis <- lm(logFev ~ logHt)
plot(xAxis$residuals,yAxis$residuals,xlab="Age residuals",
      ylab="FEV residuals")
xAxis <- lm(logHt ~ logAge); yAxis <- lm(logFev ~ logAge)
plot(xAxis$residuals,yAxis$residuals,xlab="Height residuals",
      ylab="FEV residuals")
```



## Your mission: Explore a house-price dataset

For 26 houses sold in Chicago, a long time ago, we know the selling price  $Y$  as well as eight characteristics of each house: square footage, parking information, etc.



Reference: Ashish Sen and Muni Srivastava, *Regression Analysis: Theory, Methods and Applications*, 2013.

## A peek at the house-price data

```
Q <- read.csv("houses.txt", sep=""); print(Q)
```

	##	Y	bdr	flr	fp	rms	st	lot	tax	bth	gar
## 1	53	2	967	0	5	0	39	652	1.5	0.0	
## 2	55	2	815	1	5	0	33	1000	1.0	2.0	
## 3	56	3	900	0	5	1	35	897	1.5	1.0	
## 4	58	3	1007	0	6	1	24	964	1.5	2.0	
## 5	64	3	1100	1	7	0	50	1099	1.5	1.5	
## 6	44	4	897	0	7	0	25	960	2.0	1.0	
## 7	49	5	1400	0	8	0	30	678	1.0	1.0	
## 8	70	3	2261	0	6	0	29	2700	1.0	2.0	
## 9	72	4	1290	0	8	1	33	800	1.5	1.5	
## 10	82	4	2104	0	9	0	40	1038	2.5	1.0	
## 11	85	8	2240	1	12	1	50	1200	3.0	2.0	
## 12	45	2	641	0	5	0	25	860	1.0	0.0	
## 13	47	3	862	0	6	0	25	600	1.0	0.0	
## 14	49	4	1043	0	7	0	30	676	1.5	0.0	
## 15	56	4	1325	0	8	0	50	1287	1.5	0.0	
## 16	60	2	782	0	5	1	25	834	1.0	0.0	
## 17	62	3	1126	0	7	1	30	734	2.0	0.0	
## 18	64	4	1226	0	8	0	37	551	2.0	2.0	
## 19	66	2	929	1	5	0	30	1355	1.0	1.0	
## 20	35	4	1137	0	7	0	25	561	1.5	0.0	

## The response variable and eight predictors

Variable	Meaning
Y	Selling price in thousands of dollars
bdr	Number of bedrooms
flr	Floor space in square feet
fp	Number of fireplaces
rms	Number of rooms
st	Storm windows present (indicator variable)
lot	Frontage in feet
bth	Number of bathrooms
gar	Number of garage parking spaces

Next week, we'll use the techniques on these slides to analyse the dataset. For those wanting a head start, the data are available on Portal.

# Course evaluations

Please visit <https://courseevaluations.utoronto.ca>



## U of T Course Evaluations

Centre for Teaching Support & Innovation – 130 St. George Street, Robarts Library, 4th floor

## The Course Evaluation Framework at the University of Toronto

Welcome to U of T Course Evaluations!

**Students: Click here to complete your course evaluations on PORTAL!**

The University of Toronto is committed to ensuring the quality of its academic programs, its teaching, and the learning experiences of its students. An essential component of our commitment to teaching excellence is the regular evaluation of courses by students. At the University of Toronto, course evaluations are conducted to collect formative data for instructors to improve their teaching, to provide summative data for administrative purposes (such as annual merit, tenure, and promotion review) and for program and curriculum review, and to provide members of the University community, including students, with information about teaching and courses at the university.

## Next steps

- ▶ Try questions 1 to 5 in Chapter 6, using the techniques from our lecture slides. Solutions will be posted by Wednesday 29 November
- ▶ The exam on 19 December will contain 25% pre-midterm content: questions which you could answer without having attended lectures after mid-October or having read §3.3 onwards



## Postscript: Assessing linearity visually in the meadowfoam dataset

The following R code accompanies our thinking on slide 41:

```
Flowers0<-case0901$Flowers[case0901$Time==1]
Flowers1<-case0901$Flowers[case0901$Time==2]
Intensity0<-case0901$Intensity[case0901$Time==1]
Intensity1<-case0901$Intensity[case0901$Time==2]
plot (Intensity1,Flowers1,ylab="Flowers",xlab="Intensity",
      col="green",pch=19,ylim=c(31,79))
lines (Intensity0,Flowers0,ylab="Flowers",xlab="Intensity",
       col="red",type="p",pch=19)
legend(700,75,c("Early timing","Late timing"),pch=c(19,19),
       col=c("green","red"))
```

## Postscript: Assessing linearity visually in the meadowfoam dataset

