

# Multiple Linear Regression (MLR) Handouts

Yibi Huang

- Data and Models
- Least Squares Estimate, Fitted Values, Residuals
- Sum of Squares
- Do Regression in R
- Interpretation of Regression Coefficients
- $t$ -Tests on Individual Regression Coefficients
- $F$ -Tests on Multiple Regression Coefficients/Goodness-of-Fit

## Data for Multiple Linear Regression

Multiple linear regression is a generalized form of simple linear regression, in which the data contains multiple explanatory variables.

	SLR		MLR				
	$\mathbf{x}$	$\mathbf{y}$	$\mathbf{x}_1$	$\mathbf{x}_2$	$\dots$	$\mathbf{x}_p$	$\mathbf{y}$
case 1:	$x_1$	$y_1$	$x_{11}$	$x_{12}$	$\dots$	$x_{1p}$	$y_1$
case 2:	$x_2$	$y_2$	$x_{21}$	$x_{22}$	$\dots$	$x_{2p}$	$y_2$
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
case $n$ :	$x_n$	$y_n$	$x_{n1}$	$x_{n2}$	$\dots$	$x_{np}$	$y_n$

- ▶ For SLR, we observe pairs of variables.  
For MLR, we observe rows of variables.  
Each row (or pair) is called a *case*, a *record*, or a *data point*
- ▶  $y_i$  is the *response* (or *dependent variable*) of the  $i$ th observation
- ▶ There are  $p$  *explanatory variables* (or *covariates*, *predictors*, *independent variables*), and  $x_{ik}$  is the value of the explanatory variable  $\mathbf{x}_k$  of the  $i$ th case

# Multiple Linear Regression Models

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i \quad \text{where } \varepsilon_i \text{'s are i.i.d. } N(0, \sigma^2)$$

In the model above,

- ▶  $\varepsilon_i$ 's (errors, or noise) are i.i.d.  $N(0, \sigma^2)$
- ▶ Parameters include:

$\beta_0 =$  intercept;

$\beta_k =$  regression coefficients (slope) for the  $k$ th explanatory variable,  $k = 1, \dots, p$

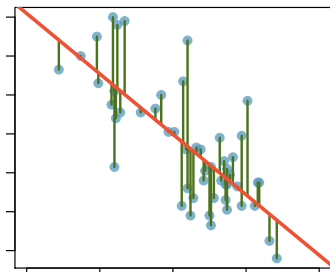
$\sigma^2 = \text{Var}(\varepsilon_i)$  is the variance of errors

- ▶ Observed (known):  $y_i, x_{i1}, x_{i2}, \dots, x_{ip}$   
Unknown:  $\beta_0, \beta_1, \dots, \beta_p, \sigma^2, \varepsilon_i$ 's
- ▶ Random variables:  $\varepsilon_i$ 's,  $y_i$ 's  
Constants (nonrandom):  $\beta_k$ 's,  $\sigma^2, x_{ik}$ 's

## Fitting the Model — Least Squares Method

Recall for SLR, the least squares estimate  $(\hat{\beta}_0, \hat{\beta}_1)$  for  $(\beta_0, \beta_1)$  is the intercept and slope of the straight line with the minimum sum of squared vertical distance to the data points

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2.$$



MLR is just like SLR. The least squares estimate  $(\hat{\beta}_0, \dots, \hat{\beta}_p)$  for  $(\beta_0, \dots, \beta_p)$  is the intercept and slopes of the (hyper)plane with the minimum sum of squared vertical distance to the data points

$$\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip})^2$$

## Solving the Least Squares Problem (1)

From now on, we use the “hat” symbol to differentiate the estimated coefficient  $\hat{\beta}_j$  from the actual unknown coefficient  $\beta_j$ .

---

To find the  $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$  that minimize

$$L(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p) = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip})^2$$

one can set the derivatives of  $L$  with respect to  $\hat{\beta}_j$  to 0

$$\frac{\partial L}{\partial \hat{\beta}_0} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip})$$

$$\frac{\partial L}{\partial \hat{\beta}_k} = -2 \sum_{i=1}^n x_{ik} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip}), \quad k = 1, 2, \dots, p$$

and then equate them to 0. This results in a system of  $(p + 1)$  equations in  $(p + 1)$  unknowns.

## Solving the Least Squares Problem (2)

The least squares estimate  $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$  is the solution to the following system of equations, called *normal equations*.

$$\begin{aligned} n\hat{\beta}_0 &+ \hat{\beta}_1 \sum_{i=1}^n x_{i1} &+ \dots &+ \hat{\beta}_p \sum_{i=1}^n x_{ip} &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_{i1} &+ \hat{\beta}_1 \sum_{i=1}^n x_{i1}^2 &+ \dots &+ \hat{\beta}_p \sum_{i=1}^n x_{i1}x_{ip} &= \sum_{i=1}^n x_{i1}y_i \\ &&&\vdots & \\ \hat{\beta}_0 \sum_{i=1}^n x_{ik} &+ \hat{\beta}_1 \sum_{i=1}^n x_{ik}x_{i1} &+ \dots &+ \hat{\beta}_p \sum_{i=1}^n x_{ik}x_{ip} &= \sum_{i=1}^n x_{ik}y_i \\ &&&\vdots & \\ \hat{\beta}_0 \sum_{i=1}^n x_{ip} &+ \hat{\beta}_1 \sum_{i=1}^n x_{ip}x_{i1} &+ \dots &+ \hat{\beta}_p \sum_{i=1}^n x_{ip}^2 &= \sum_{i=1}^n x_{ip}y_i \end{aligned}$$

- ▶ Don't worry about solving the equations.  
R and many other softwares can do the computation for us.
- ▶ In general,  $\hat{\beta}_j \neq \beta_j$ , but they will be close under some conditions

## Fitted Values

The fitted value or predicted value:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip}$$

- ▶ Again, the “hat” symbol is used to differentiate the fitted value  $\hat{y}_i$  from the actual observed value  $y_i$ .

# Residuals

- ▶ One cannot directly compute the errors

$$\varepsilon_i = y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip}$$

since the coefficients  $\beta_0, \beta_1, \dots, \beta_p$  are *unknown*.

- ▶ The errors  $\varepsilon_i$  can be estimated by the **residuals**  $e_i$  defined as:

$$\begin{aligned} \text{residual } e_i &= \text{observed } y_i - \text{predicted } y_i \\ &= y_i - \hat{y}_i \\ &= y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip}) \\ &= \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i \\ &\quad - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip} \end{aligned}$$

- ▶  $e_i \neq \varepsilon_i$  in general since  $\hat{\beta}_j \neq \beta_j$
- ▶ Graphical explanation



# Properties of Residuals

Recall the least squares estimate  $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$  satisfies the equations

$$\sum_{i=1}^n \underbrace{(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip})}_{= y_i - \hat{y}_i = e_i = \text{residual}} = 0 \text{ and}$$

$$\sum_{i=1}^n x_{ik} \underbrace{(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip})}_{= y_i - \hat{y}_i = e_i = \text{residual}} = 0, \quad k = 1, 2, \dots, p.$$

Thus the residuals  $e_i$  have the properties

$$\underbrace{\sum_{i=1}^n e_i}_{\text{Residuals add up to 0.}} = 0, \quad \underbrace{\sum_{i=1}^n x_{ik} e_i}_{\text{Residuals are orthogonal to covariates.}} = 0, \quad k = 1, 2, \dots, p.$$

## Sum of Squares

Observe that

$$y_i - \bar{y} = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i)$$

Squaring up both sides we get

$$(y_i - \bar{y})^2 = (\hat{y}_i - \bar{y})^2 + (y_i - \hat{y}_i)^2 + 2(\hat{y}_i - \bar{y})(y_i - \hat{y}_i)$$

Summing up over all the cases  $i = 1, 2, \dots, n$ , we get

$$\begin{aligned} \overbrace{\sum_{i=1}^n (y_i - \bar{y})^2}^{\text{SST}} &= \overbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}^{\text{SSR}} + \overbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}^{\text{SSE}} \\ &\quad + 2 \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})(y_i - \hat{y}_i)}_{= 0, \text{ see next page.}} \end{aligned}$$

Why  $\sum_{i=1}^n (\hat{y}_i - \bar{y})(y_i - \hat{y}_i) = 0$ ?

$$\begin{aligned} & \sum_{i=1}^n (\hat{y}_i - \bar{y}) \underbrace{(y_i - \hat{y}_i)}_{=e_i} \\ &= \sum_{i=1}^n \hat{y}_i e_i - \sum_{i=1}^n \bar{y} e_i \\ &= \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip}) e_i - \sum_{i=1}^n \bar{y} e_i \\ &= \hat{\beta}_0 \underbrace{\sum_{i=1}^n e_i}_{=0} + \hat{\beta}_1 \underbrace{\sum_{i=1}^n x_{i1} e_i}_{=0} + \dots + \hat{\beta}_p \underbrace{\sum_{i=1}^n x_{ip} e_i}_{=0} - \bar{y} \underbrace{\sum_{i=1}^n e_i}_{=0} \\ &= 0 \end{aligned}$$

in which we used the properties of residuals that  $\sum_{i=1}^n e_i = 0$  and  $\sum_{i=1}^n x_{ik} e_i = 0$  for all  $k = 1, \dots, p$ .

## Interpretation of Sum of Squares

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{SST}} = \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\text{SSR}} + \underbrace{\sum_{i=1}^n \overbrace{(y_i - \hat{y}_i)}^{=e_i}}_{\text{SSE}}$$

- ▶ SST = **total sum of squares**
  - ▶ total variability of  $\mathbf{y}$
  - ▶ depends on the response  $\mathbf{y}$  only, not on the form of the model
- ▶ SSR = **regression sum of squares**
  - ▶ variability of  $\mathbf{y}$  explained by  $\mathbf{x}_1, \dots, \mathbf{x}_p$
- ▶ SSE = **error (residual) sum of squares**
  - ▶  $= \min_{\beta_0, \beta_1, \dots, \beta_p} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2$
  - ▶ variability of  $\mathbf{y}$  not explained by  $\mathbf{x}$ 's

## Degrees of Freedom

If the MLR model  $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i$ ,  $\varepsilon_i$ 's i.i.d.  $\sim N(0, \sigma^2)$  is true, it can be shown that

$$\frac{\text{SSE}}{\sigma^2} \sim \chi_{n-p-1}^2,$$

If we further assume that  $\beta_1 = \beta_2 = \dots = \beta_p = 0$ , then

$$\frac{\text{SST}}{\sigma^2} \sim \chi_{n-1}^2, \quad \frac{\text{SSR}}{\sigma^2} \sim \chi_p^2$$

and SSR is independent of SSE.

Note the **degrees of freedom** of the 3 chi-square distributions

$$dfT = n - 1, \quad dfR = p, \quad dfE = n - p - 1$$

break down similarly

$$dfT = dfR + dfE$$

just like  $\text{SST} = \text{SSR} + \text{SSE}$ .

## Why SSE Has $n - p - 1$ Degrees of Freedom?

The  $n$  residuals  $e_1, \dots, e_n$  cannot all vary freely.

There are  $p + 1$  constraints:

$$\sum_{i=1}^n e_i = 0 \quad \text{and} \quad \sum_{i=1}^n x_{ki} e_i = 0 \quad \text{for } k = 1, \dots, p.$$

So only  $n - (p + 1)$  of them can be *freely varying*.

The  $p + 1$  constraints comes from the  $p + 1$  coefficients  $\beta_0, \dots, \beta_p$  in the model, and each contributes one constraint  $\frac{\partial}{\partial \beta_k} = 0$ .

## Mean Square Error (MSE) — Estimate of $\sigma^2$

The **mean squares** is the sum of squares divided by its degrees of freedom:

$$MST = \frac{SST}{dfT} = \frac{SST}{n-1} = \text{sample variance of } Y,$$

$$MSR = \frac{SSR}{dfR} = \frac{SSR}{p},$$

$$MSE = \frac{SSE}{dfE} = \frac{SSE}{n-p-1} = \hat{\sigma}^2$$

- ▶ From the fact  $\frac{SSE}{\sigma^2} \sim \chi_{n-p-1}^2$  and that the mean of a  $\chi_k^2$  distribution is  $k$ , we know that **MSE is an unbiased estimator for  $\sigma^2$** .
- ▶ Though SSE always decreases as we add terms to the model, adding unimportant terms may increase MSE.

## Example: Housing Price

Price	BDR	FLR	FP	RMS	ST	LOT	BTH	CON	GAR	LOC	
53	2	967	0	5	0	39	1.5	1	0.0	0	
55	2	815	1	5	0	33	1.0	1	2.0	0	Price = Selling price in \$1000
56	3	900	0	5	1	35	1.5	1	1.0	0	BDR = Number of bedrooms
58	3	1007	0	6	1	24	1.5	0	2.0	0	FLR = Floor space in sq. ft.
64	3	1100	1	7	0	50	1.5	1	1.5	0	FP = Number of fireplaces
44	4	897	0	7	0	25	2.0	0	1.0	0	RMS = Number of rooms
49	5	1400	0	8	0	30	1.0	0	1.0	0	ST = Storm windows
70	3	2261	0	6	0	29	1.0	0	2.0	0	(1 if present, 0 if absent)
72	4	1290	0	8	1	33	1.5	1	1.5	0	
82	4	2104	0	9	0	40	2.5	1	1.0	0	LOT = Front footage of lot in feet
85	8	2240	1	12	1	50	3.0	0	2.0	0	BTH = Number of bathrooms
45	2	641	0	5	0	25	1.0	0	0.0	1	CON = Construction
47	3	862	0	6	0	25	1.0	1	0.0	1	(1 if frame, 0 if brick)
49	4	1043	0	7	0	30	1.5	0	0.0	1	GAR = Garage size
56	4	1325	0	8	0	50	1.5	0	0.0	1	(0 = no garage,
60	2	782	0	5	1	25	1.0	0	0.0	1	1 = one-car garage, etc.)
62	3	1126	0	7	1	30	2.0	1	0.0	1	
64	4	1226	0	8	0	37	2.0	0	2.0	1	LOC = Location
.											(1 if property is in zone A,
.											0 otherwise)
.											
50	2	691	0	6	0	30	1.0	0	2.0	0	
65	3	1023	0	7	1	30	2.0	1	1.0	0	



# How to Do Regression Using R?

```
> housing = read.table("housing.txt",h=TRUE) # to load the data  
> lm(Price ~ FLR+LOT+BDR+GAR+ST, data=housing)
```

Call:

```
lm(formula = Price ~ FLR + LOT + BDR + GAR + ST, data = housing)
```

Coefficients:

(Intercept)	FLR	LOT	BDR	GAR	ST
24.63232	0.02009	0.44216	-3.44509	3.35274	11.64033

---

The `lm()` command above asks R to fit the model

$$\text{Price} = \beta_0 + \beta_1\text{FLR} + \beta_2\text{LOT} + \beta_3\text{BDR} + \beta_4\text{GAR} + \beta_5\text{ST} + \varepsilon$$

and R gives us the regression equation

$$\widehat{\text{Price}} = 24.63 + 0.02\text{FLR} + 0.44\text{LOT} - 3.45\text{BDR} + 3.35\text{GAR} + 11.64\text{ST}$$

$$\widehat{\text{Price}} = 24.63 + 0.02\text{FLR} + 0.44\text{LOT} - 3.45\text{BDR} + 3.35\text{GAR} + 11.64\text{ST}$$

The regression equation tells us:

- ▶ an extra square foot in floor area increases the price by \$20,
- ▶ an extra foot in front footage by ..... \$440,
- ▶ an additional bedroom by ..... -\$3450,
- ▶ an additional space in the garage by ..... \$3350,
- ▶ using storm windows by ..... \$11640.

**Question:**

Why an additional bedroom makes a house less valuable?

# Interpretation of Regression Coefficients

- ▶  $\beta_0$  = intercept = the mean value of  $y$  when all  $x_j$ ' are 0.
  - ▶ may not have practical meaning  
e.g.,  $\beta_0$  is meaningless in the housing price model as no housing unit has 0 floor space.
- ▶  $\beta_j$ : regression coefficient for  $x_j$ , is the mean change in the response  $y$  when  $x_j$  is increased by one unit *holding all other  $x_j$ 's constant*
  - ▶ Interpretation of  $\beta_j$  depends on the presence of other covariates in the model  
e.g., the meaning of the 2  $\beta_1$ 's in the following 2 models are different

$$\text{Model 1 : } y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i$$

$$\text{Model 2 : } y_i = \beta_0 + \beta_1 x_{i1} + \varepsilon_i.$$

# What's Wrong?

```
# Model 1  
> lm(Price ~ BDR, data=housing)
```

(Intercept)	BDR
43.487	3.921

The regression coefficient for BDR is 3.921 in the Model 1 above but  $-3.445$  in the Model 2 below.

```
# Model 2  
> lm(Price ~ FLR+LOT+BDR+GAR+ST, data=housing)
```

(Intercept)	FLR	LOT	BDR	GAR	ST
24.63232	0.02009	0.44216	-3.44509	3.35274	11.64033

Considering **BDR** alone, house prices *increase* with **BDR**.

However, an extra bedroom makes a housing unit less valuable when when other covariates (**FLR**, **LOT**, etc) are fixed.

Does this make sense?

## More R Commands

```
> lm1 = lm(Price ~ FLR+RMS+BDR+GAR+LOT+ST+CON+LOC, data=housing)
> summary(lm1)      # Regression output with more details
                    # including multiple R-squared,
                    # and the estimate of sigma

> lm1$coef          # show the estimated beta's
> lm1$fitted        # show the fitted values
> lm1$res           # show the residuals
```

```
> lm1 = lm(Price ~ FLR+LOT+BDR+GAR+ST, data=housing)
> summary(lm1)
```

Call:

```
lm(formula = Price ~ FLR + LOT + BDR + GAR + ST, data = housing)
```

Residuals:

Min	1Q	Median	3Q	Max
-9.7530	-2.9535	0.1779	3.7183	12.9728

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	24.632318	4.836743	5.093	5.56e-05	***
FLR	0.020094	0.003668	5.478	2.31e-05	***
LOT	0.442164	0.150023	2.947	0.007965	**
BDR	-3.445086	1.279347	-2.693	0.013995	*
GAR	3.352739	1.560239	2.149	0.044071	*
ST	11.640334	2.688867	4.329	0.000326	***

---

Residual standard error: 5.79 on 20 degrees of freedom  
Multiple R-squared: 0.8306, Adjusted R-squared: 0.7882  
F-statistic: 19.61 on 5 and 20 DF, p-value: 4.306e-07

## *t*-Tests on Individual Regression Coefficients

For a MLR model  $Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \varepsilon_i$ , to test the hypotheses,

$$H_0 : \beta_j = c \quad \text{v.s.} \quad H_a : \beta_j \neq c$$

the *t*-statistic is

$$t = \frac{\widehat{\beta}_j - c}{\text{SE}(\widehat{\beta}_j)}$$

in which  $\text{SE}(\widehat{\beta}_j)$  is the standard error for  $\widehat{\beta}_j$ .

- ▶ General formula for  $\text{SE}(\widehat{\beta}_j)$  is a bit complicated but unimportant in STAT222 and hence is omitted
- ▶ R can compute  $\text{SE}(\widehat{\beta}_j)$  for us
- ▶ Formula for  $\text{SE}(\widehat{\beta}_j)$  for a few special cases will be given later

This *t*-statistic also has a *t*-distribution with  $n - p - 1$  degrees of freedom

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	24.632318	4.836743	5.093	5.56e-05	***
FLR	0.020094	0.003668	5.478	2.31e-05	***
LOT	0.442164	0.150023	2.947	0.007965	**
.....(some rows are omitted)					
ST	11.640334	2.688867	4.329	0.000326	***

- ▶ the first column gives variable names
- ▶ the column `Estimate` gives the LS estimate  $\hat{\beta}_j$ 's for  $\beta_j$ 's
- ▶ the column `Std. Error` gives  $SE(\hat{\beta}_j)$ , the standard error of  $\hat{\beta}_j$
- ▶ the column `t value` gives  $t\text{-value} = \frac{\hat{\beta}_j}{SE(\hat{\beta}_j)}$
- ▶ column `Pr(>|t|)` gives the  $P$ -value for testing  $H_0: \beta_j = 0$  v.s.  $H_a: \beta_j \neq 0$ .

E.g., for `LOT`, we see

$$\hat{\beta}_{LOT} \approx 0.442, SE(\hat{\beta}_{LOT}) \approx 0.150, t = \frac{\hat{\beta}_{LOT}}{SE(\hat{\beta}_{LOT})} \approx \frac{0.442}{0.150} \approx 2.947.$$

The  $P$ -value 0.007965 is the 2-sided  $P$ -value for testing  $H_0: \beta_{LOT} = 0$



## Nested Models

We say Model 1 is **nested in** Model 2 if Model 1 is a special case of Model 2 (and hence Model 2 is an extension of Model 1).

E.g., for the 4 models below,

$$\text{Model A : } Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \varepsilon$$

$$\text{Model B : } Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

$$\text{Model C : } Y = \beta_0 + \beta_1 X_1 + \beta_3 X_3 + \varepsilon$$

$$\text{Model D : } Y = \beta_0 + \beta_1 (X_1 + X_2) + \varepsilon$$

- ▶ B is nested in A . . . . . since A reduces to B when  $\beta_3 = 0$
- ▶ C is also nested in A . . . . . since A reduces to C when  $\beta_2 = 0$
- ▶ D is nested in B . . . . . since B reduces to D when  $\beta_1 = \beta_2$
- ▶ B and C are NOT nested in either way
- ▶ D is NOT nested in C

## Nested Relationship is Transitive

If Model 1 is nested in Model 2, and Model 2 is nested in Model 3, then Model 1 is also nested in Model 3.

For example, for models in the previous slide,

D is nested in B, and B is nested in A,

implies D is also nested in A, which is clearly true because Model A reduces to Model D when

$$\beta_1 = \beta_2, \text{ and } \beta_3 = 0.$$

---

When two models are nested (Model 1 is nested in Model 2),

- ▶ the smaller model (Model 1) is called the **reduced model**, and
- ▶ the more general model (Model 2) is called the **full model**.

## SST of Nested Models

**Question:** Compare the SST for Model A and the SST for Model B. Which one is larger? Or are they equal?

What about the SST for Model C? For Model D?

## SSE of Nested Models

When a reduced model is nested in a full model, then

$$(i) SSE_{reduced} \geq SSE_{full}, \text{ and } (ii) SSR_{reduced} \leq SSR_{full}.$$

*Proof.* We will prove (i) for

- ▶ the full model  $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \varepsilon_i$  and
- ▶ the reduced model  $y_i = \beta_0 + \beta_1 x_{i1} + \beta_3 x_{i3} + \varepsilon_i$ .

The proofs for other nested models are similar.

$$\begin{aligned} SSE_{full} &= \min_{\beta_0, \beta_1, \beta_2, \beta_3} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \beta_3 x_{i3})^2 \\ &\leq \min_{\beta_0, \beta_1, \beta_3} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \beta_3 x_{i3})^2 \\ &= SSE_{reduced} \end{aligned}$$

Part (ii) follows directly from (i), the identity  $SST = SSR + SSE$ , and the fact that all MLR models of the same data set have a common  $SST$

## General Framework for Testing Nested Models

$H_0$ : reduced model is true v.s.  $H_a$ : full model is true

- ▶ As the reduced model is nested in the full model,

$$SSE_{reduced} \geq SSE_{full}$$

- ▶ Simplicity or Accuracy?
  - ▶ The full model fits the data better (with smaller SSE) but is more complicate
  - ▶ The reduced model doesn't fit as well but is simpler.
  - ▶ If  $SSE_{reduced} \approx SSE_{full}$ , one can sacrifice a bit of accuracy in exchange for simplicity
  - ▶ If  $SSE_{reduced} \gg SSE_{full}$ , it would cost to much in accuracy in exchange for simplicity. The full model is preferred.

## The $F$ -Statistic

$$F = \frac{(SSE_{reduced} - SSE_{full}) / (df_{reduced} - df_{full})}{MSE_{full}}$$

- ▶  $SSE_{reduced} - SSE_{full}$  is the reduction in SSE from replacing the reduced model with the full model.
- ▶  $df_{reduced}$  is the df for **error** for the reduced model.
- ▶  $df_{full}$  is the df for **error** for the full model.
- ▶  $F \geq 0$  since  $SSE_{reduced} \geq SSE_{full} \geq 0$
- ▶ The smaller the  $F$ -statistic, the more we favor the reduced model
- ▶ Under  $H_0$ , the  $F$ -statistic has an  $F$ -distribution with  $df_{reduced} - df_{full}$  and  $df_{full}$  degrees of freedom.

## Testing All Coefficients Equal Zero

Testing the hypotheses

$$H_0: \beta_1 = \cdots = \beta_p = 0 \text{ v.s. } H_a: \text{ not all } \beta_1, \dots, \beta_p = 0$$

is a test to evaluate the **overall significance** of a model.

$$\text{Full : } y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \varepsilon_i$$

$$\text{Reduced : } y_i = \beta_0 + \varepsilon_i \quad (\text{all covariates are unnecessary})$$

- ▶ The LS estimate for  $\beta_0$  in the reduced model is  $\hat{\beta}_0 = \bar{y}$ , so

$$SSE_{\text{reduced}} = \sum_{i=1}^n (y_i - \hat{\beta}_0)^2 = \sum_i (y_i - \bar{y})^2 = SST_{\text{full}}$$

- ▶  $df_{\text{reduced}} = dfE_{\text{reduced}} = n - 1$ ,  
because the reduced model has only one coefficient  $\beta_0$
- ▶  $df_{\text{full}} = dfE_{\text{full}} = n - p - 1$ .

## Testing All Coefficients Equal Zero

Hence

$$\begin{aligned} F &= \frac{(SSE_{reduced} - SSE_{full}) / (df_{reduced} - df_{full})}{MSE_{full}} \\ &= \frac{(SST_{full} - SSE_{full}) / [n - 1 - (n - p - 1)]}{SSE_{full} / (n - p - 1)} \\ &= \frac{SSR_{full} / p}{SSE_{full} / (n - p - 1)}. \end{aligned}$$

Moreover,  $F \sim F_{p, n-p-1}$  under  $H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$ .

In R, the  $F$  statistic and  $p$ -value are displayed in the last line of the output of the `summary()` command.

```
> lm1 = lm(Price ~ FLR+LOT+BDR+GAR+ST, data=housing)
> summary(lm1)
... (output omitted)
```

```
Residual standard error: 5.79 on 20 degrees of freedom
Multiple R-squared: 0.8306, Adjusted R-squared: 0.7882
F-statistic: 19.61 on 5 and 20 DF, p-value: 4.306e-07
```



## ANOVA and the $F$ -Test

The test of all coefficients equal zero is often summarized in an ANOVA table.

Source	df	Sum of Squares	Mean Squares	$F$
Regression	$dfR = p$	SSR	$MSR = \frac{SSR}{dfR}$	$F = \frac{MSR}{MSE}$
Error	$dfE = n - p - 1$	SSE	$MSE = \frac{SSE}{dfE}$	
Total	$dfT = n - 1$	SST		

## Testing Some Coefficients Equal to Zero

E.g., for the housing price data, we may want to test if we can eliminate **BDR** and **GAR** from the model,

i.e.,  $H_0: \beta_{BDR} = \beta_{GAR} = 0$ .

```
> lmfull = lm(Price ~ FLR+LOT+BDR+GAR+ST, data=housing)
> lmreduced = lm(Price ~ FLR+LOT+ST, data=housing)
> anova(lmreduced, lmfull)
Analysis of Variance Table
```

```
Model 1: Price ~ FLR + LOT + ST
```

```
Model 2: Price ~ FLR + LOT + BDR + GAR + ST
```

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	22	1105.01				
2	20	670.55	2	434.46	6.4792	0.006771 **

Note SSE is called RSS (residual sum of square) in R.

## Testing Equality of Coefficients

Example. To test  $H_0: \beta_1 = \beta_2 = \beta_3$ , the reduced model is

$$\begin{aligned} Y &= \beta_0 + \beta_1 X_1 + \beta_1 X_2 + \beta_1 X_3 + \beta_4 X_4 + \varepsilon \\ &= \beta_0 + \beta_1 (X_1 + X_2 + X_3) + \beta_4 X_4 + \varepsilon \end{aligned}$$

1. Make a new variable  $W = X_1 + X_2 + X_3$
2. Fit the reduced model by regressing  $Y$  on  $W$  and  $X_4$
3. Find  $SSE_{reduced}$  and  $df_{reduced} - df_{full} = \underline{2}$
4. In R

```
> lmfull = lm(Y ~ X1 + X2 + X3 + X4)
> lmreduced = lm(Y ~ I(X1 + X2 + X3) + X4)
> anova(lmreduced, lmfull)
```

The line `lmreduced = lm(Y ~ I(X1 + X2 + X3) + X4)` is equivalent to

```
> W = X1 + X2 + X3
> lmreduced = lm(Y ~ W + X4)
```