

In Chapter 15:

- 15.1 The General Idea
- 15.2 The Multiple Regression Model
- 15.3 Categorical Explanatory Variables
- 15.4 Regression Coefficients
- [15.5 ANOVA for Multiple Linear Regression]
- [15.6 Examining Conditions]

[Not covered in recorded presentation]

Basic Biostat 15: Multiple Linear Regression 2
© 2008 Jones and Bartlett Publishers

15.1 The General Idea

Simple regression considers the relation between a single explanatory variable and response variable

$X \rightarrow Y$

Basic Biostat 15: Multiple Linear Regression 3
© 2008 Jones and Bartlett Publishers

The General Idea

Multiple regression simultaneously considers the influence of multiple explanatory variables on a response variable Y

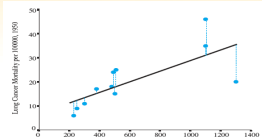
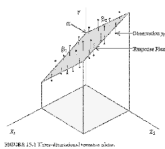
$$\begin{matrix} X_1 & \rightarrow & \\ X_2 & \rightarrow & Y \\ \vdots & & \\ X_k & \rightarrow & \end{matrix}$$

The intent is to look at the independent effect of each variable while “adjusting out” the influence of potential confounders

Basic Biostat 15: Multiple Linear Regression 4
© 2008 Jones and Bartlett Publishers

Regression Modeling

- A simple regression model (one independent variable) fits a regression *line* in 2-dimensional space
- A multiple regression model with two explanatory variables fits a regression *plane* in 3-dimensional space

Basic Biostat 15: Multiple Linear Regression 5
© 2008 Jones and Bartlett Publishers

Simple Regression Model

Regression coefficients are estimated by minimizing $\sum \text{residuals}^2$ (i.e., sum of the squared residuals) to derive this model:

$$\hat{y} = a + bx$$

The **standard error of the regression** ($s_{Y|x}$) is based on the squared residuals:

$$s_{Y|x} = \sqrt{\sum \text{residuals}^2 / df_{\text{res}}}$$

Basic Biostat 15: Multiple Linear Regression 6
© 2008 Jones and Bartlett Publishers

Multiple Regression Model

Again, **estimates for the multiple slope coefficients** are derived by minimizing $\sum \text{residuals}^2$ to derive this multiple regression model:

$$\hat{y} = a + b_1x_1 + b_2x_2$$

Again, the **standard error of the regression** is based on the $\sum \text{residuals}^2$:

$$S_{Y|x} = \sqrt{\sum \text{residuals}^2 / df_{\text{res}}}$$

Basic Biostat 15: Multiple Linear Regression 7 © 2008 Jones and Bartlett Publishers

Multiple Regression Model

- Intercept α predicts where the regression plane crosses the Y axis
- Slope for variable X_1 (β_1) predicts the change in Y per unit X_1 holding X_2 constant
- The slope for variable X_2 (β_2) predicts the change in Y per unit X_2 holding X_1 constant

Basic Biostat 15: Multiple Linear Regression 15: FIGURE 15.1 Three-dimensional response plane. © 2008 Jones and Bartlett Publishers

Multiple Regression Model

A multiple regression model with k explanatory variables fits a regression "surface" in $k + 1$ dimensional space (cannot be visualized)

Basic Biostat 15: Multiple Linear Regression 9 © 2008 Jones and Bartlett Publishers

15.3 Categorical Explanatory Variables in Regression Models

- Categorical explanatory variables can be incorporated into a regression model by converting them into 0/1 ("dummy") variables
- For binary variables, code dummies "0" for "no" and 1 for "yes"

Basic Biostat 15: Multiple Linear Regression 10 © 2008 Jones and Bartlett Publishers

Dummy Variables, More than two levels

For categorical variables with k categories, use $k-1$ dummy variables

SMOKE2 has three levels, initially coded
 0 = non-smoker
 1 = former smoker
 2 = current smoker

Use $k - 1 = 3 - 1 = 2$ dummy variables to code this information like this:

SMOKE2	DUMMY1	DUMMY2
0	0	0
1	1	0
2	0	1

Basic Biostat 15: Multiple Linear Regression 11 © 2008 Jones and Bartlett Publishers

Illustrative Example

Childhood respiratory health survey.

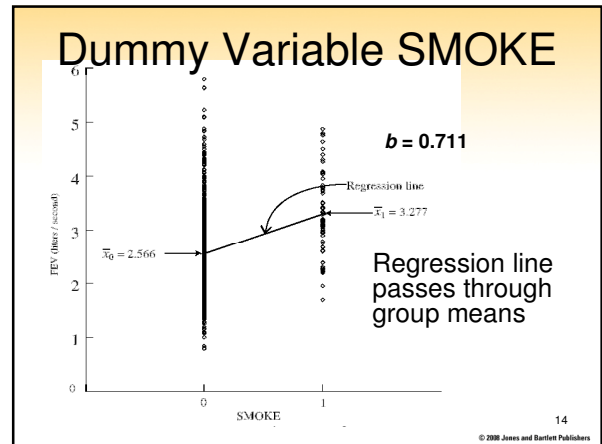
- Binary explanatory variable (SMOKE) is coded 0 for non-smoker and 1 for smoker
- Response variable Forced Expiratory Volume (FEV) is measured in liters/second
- The mean FEV in nonsmokers is 2.566
- The mean FEV in smokers is 3.277

Basic Biostat 15: Multiple Linear Regression 12 © 2008 Jones and Bartlett Publishers

Example, cont.

- Regress FEV on SMOKE least squares regression line:
 $\hat{y} = 2.566 + 0.711X$
- Intercept (2.566) = the mean FEV of group 0
- Slope = the mean difference in FEV
 $= 3.277 - 2.566 = 0.711$
- $t_{stat} = 6.464$ with 652 *df*, $P \approx 0.000$ (same as equal variance *t* test)
- The 95% CI for slope is 0.495 to 0.927 (same as the 95% CI for $\mu_1 - \mu_0$)

Basic Biostat 15: Multiple Linear Regression 13 © 2008 Jones and Bartlett Publishers



Smoking increases FEV?

- Children who smoked had higher mean FEV
- How can this be true given what we know about the deleterious respiratory effects of smoking?
- ANS: Smokers were older than the nonsmokers
- AGE confounded the relationship between SMOKE and FEV
- A multiple regression model can be used to adjust for AGE in this situation

Basic Biostat 15: Multiple Linear Regression 15 © 2008 Jones and Bartlett Publishers

15.4 Multiple Regression Coefficients

Rely on software to calculate multiple regression statistics

Basic Biostat 15: Multiple Linear Regression 16 © 2008 Jones and Bartlett Publishers

Example

SPSS output for our example:

		Coefficients ^a				
		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
Model		B	Std. Error	Beta		
1	(Constant)	0.367	.081		4.511	.000
	SMOKE	-.209	.081	-.072	-2.588	.010
	AGE	.231	.008	.786	28.176	.000

a. Dependent Variable: FEV

The multiple regression model is:
 $FEV = 0.367 + -.209(SMOKE) + .231(AGE)$

Basic Biostat 15: Multiple Linear Regression 17 © 2008 Jones and Bartlett Publishers

Multiple Regression Coefficients, cont.

- The slope coefficient associated for SMOKE is $-.206$, suggesting that smokers have $.206$ less FEV on average compared to non-smokers (after adjusting for age)
- The slope coefficient for AGE is $.231$, suggesting that each year of age is associated with an increase of $.231$ FEV units on average (after adjusting for SMOKE)

Basic Biostat 15: Multiple Linear Regression 18 © 2008 Jones and Bartlett Publishers

Inference About the Coefficients

Inferential statistics are calculated for each regression coefficient. For example, in testing $H_0: \beta_1 = 0$ (SMOKE coefficient controlling for AGE)

$t_{stat} = -2.588$ and $P = 0.010$

Model	Unstandardized Coefficients		Standardized Coefficients	t	Sig.	
	B	Std. Error	Beta			
1	(Constant)	.367	.081		4.511	.000
	smoke	-.209	.081	-.072	-2.588	.010
	age	.231	.008	.786	28.176	.000

a. Dependent Variable: fev

$df = n - k - 1 = 654 - 2 - 1 = 651$

Inference About the Coefficients

The 95% confidence interval for this slope of SMOKE controlling for AGE is -0.368 to -0.050.

Model	95% Confidence Interval for B	
	Lower Bound	Upper Bound
1	(Constant)	.207 .527
	smoke	-.368 -.050
	age	.215 .247

a. Dependent Variable: fev