

## AIU FOR 6520

### Statistical Research Design & Methods in Forensic Psychology

#### ● **Instructor:** Ivo Dinov,

Asst. Prof. of Statistics, Neurology, Psychology

AIU, UCLA, Winter 2003

[http://www.stat.ucla.edu/~dinov/courses\\_students.html](http://www.stat.ucla.edu/~dinov/courses_students.html)

FOR 6520, AIU, Ivo Dinov

Slide 1

#### ● Multiple Regression Analysis

FOR 6520, AIU, Ivo Dinov

Slide 2

#### Introduction

- We extend the concept of simple linear regression as we investigate a response  $y$  which is affected by several independent variables,  $x_1, x_2, x_3, \dots, x_k$ .
- Our objective is to use the information provided by the  $x_i$  to predict the value of  $y$ .

Slide 3

FOR 6520, AIU, Ivo Dinov

#### Example

- Let  $y$  be a student's college achievement, measured by his/her GPA. This might be a function of several predictor/explanatory variables:
  - $x_1$  = rank in high school class
  - $x_2$  = high school's overall rating
  - $x_3$  = high school GPA
  - $x_4$  = SAT scores
- We want to predict  $y$  using knowledge of  $x_1, x_2, x_3$  and  $x_4$ .

Slide 4

FOR 6520, AIU, Ivo Dinov

#### Example

- Let  $y$  be the monthly sales revenue for a company. This might be a function of several variables:
  - $x_1$  = advertising expenditure
  - $x_2$  = time of year
  - $x_3$  = state of economy
  - $x_4$  = size of inventory
- We want to predict  $y$  using knowledge of  $x_1, x_2, x_3$  and  $x_4$ .

Slide 5

FOR 6520, AIU, Ivo Dinov

#### Questions

- How well does the model fit?
- How strong is the relationship between the response  $y$  and the predictor variables,  $x_k$ ?
- Have any assumptions been violated?
- How good are the estimates and predictions?

We collect information using  $n$  observations on the response  $y$  and the independent variables,  $x_1, x_2, x_3, \dots, x_k$ .

Slide 6

FOR 6520, AIU, Ivo Dinov

## The General Linear Model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

- where

- ✓  $y$  is the response variable you want to predict.
- ✓  $\beta_0, \beta_1, \beta_2, \dots, \beta_k$  are unknown constants (regression parameters).
- ✓  $x_1, x_2, \dots, x_k$  are independent predictor variables, measured without error.

Slide 7

FOR6520, AUU, Ivo Dinov

## The Random Error

- The **deterministic** part of the model,
  - $E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$ ,
- describes average value of  $y$  for any fixed values of  $x_1, x_2, \dots, x_k$ . The population of measurements is generated as  $y$  deviates from the **line of means** by an amount  $\varepsilon$ . We assume
  - ✓  $\varepsilon$  are independent and identically distributed (IID)
  - ✓ Have a mean 0 and common variance  $\sigma^2$  for any set  $x_1, x_2, \dots, x_k$ .
  - ✓ Have a normal distribution.

Slide 8

FOR6520, AUU, Ivo Dinov

## Example – 2D

- Consider the model  $E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$
- This is a **first order model** (independent variables appear only to the first power).
- $\beta_0 = y\text{-intercept}$  = value of  $E(y)$  when  $x_1 = x_2 = 0$ .
- $\beta_1$  and  $\beta_2$  are the **partial regression coefficients**—the change in  $y$  for a one-unit change in  $x_i$  **when the other independent variables are held constant**.
- Traces a **plane** in three dimensional space.

Slide 9

FOR6520, AUU, Ivo Dinov

## The Method of Least Squares

- The best-fitting prediction equation is calculated using a set of  $n$  measurements ( $y, x_1, x_2, \dots, x_k$ ) as
 
$$\hat{y} = b_0 + b_1 x_1 + \dots + b_k x_k$$
- We choose our estimates  $b_0, b_1, \dots, b_k$  to estimate  $\beta_0, \beta_1, \dots, \beta_k$  to minimize

$$\begin{aligned} \text{SSE} &= \sum (y - \hat{y})^2 \\ &= \sum (y - b_0 - b_1 x_1 - \dots - b_k x_k)^2 \end{aligned}$$

Slide 10

FOR6520, AUU, Ivo Dinov

## Example – house prices

- A computer database in a small community contains the listed selling price  $y$  (in thousands of dollars), the amount of living area  $x_1$  (in hundreds of square feet), and the number of floors  $x_2$ , bedrooms  $x_3$ , and bathrooms  $x_4$ , for  $n = 15$  randomly selected residences currently on the market.

Property	$y$	$x_1$	$x_2$	$x_3$	$x_4$
1	69.0	6	1	2	1
2	118.5	10	1	2	2
3	116.5	10	1	3	2
4	130	11	2	1	2
5	90	10	1	1	1
...	...	...	...	...	...
15	209.9	21	2	4	3

Fit a first order model to the data using the method of least squares.

Use WebStat to calculate the multiple-linear regression model

Slide 11

FOR6520, AUU, Ivo Dinov

## Example

- The first order model is
 
$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4$$
 fit using WebStat with the values of  $y$  and the four independent variables entered into five columns of the Minitab worksheet.

Regression equation

### Regression Analysis: ListPrice versus SqFeet, NumFlrs, Bdrms, Baths

The regression equation is

ListPrice = 18.8 + 6.27 SqFeet - 16.2 NumFlrs - 2.67 Bdrms + 30.3 Baths

Predictor	Coef	SE Coef	T	P
Constant	18.763	9.207	2.04	0.069
SqFeet	6.2698	0.7252	8.65	0.000
NumFlrs	-16.203	0.995	-16.27	0.000
Bdrms	-2.673	0.565	-4.73	0.000
Baths	30.271	0.001	30.27	0.000

Partial regression coefficients

Slide 12

FOR6520, AUU, Ivo Dinov

## The Analysis of Variance

- The total variation in the experiment is measured by the **total sum of squares**:

$$\text{Total SS} = S_{yy} = \sum (y - \bar{y})^2$$

The **Total SS** is divided into two parts:

- ✓ **SSR** (sum of squares for regression): measures the variation explained by using the regression equation.
- ✓ **SSE** (sum of squares for error): measures the leftover variation not explained by the independent variables.

Slide 13

FOR6520, AUU, Ivo Dinov

## The ANOVA Table

$$\text{Total } df = n - 1$$

$$\text{Regression } df = k$$

$$\text{Error } df = n - 1 - k = n - k - 1$$

Mean Squares

$$\text{MSR} = \text{SSR}/k$$

$$\text{MSE} = \text{SSE}/(n-k-1)$$

Source	df	SS	MS	F
Regression	k	SSR	SSR/k	MSR/MSE
Error	n - k - 1	SSE	SSE/(n-k-1)	
Total	n - 1	Total SS		

Slide 14

FOR6520, AUU, Ivo Dinov

## The Real Estate Problem

Another portion of the SOCR printout shows the ANOVA Table, with  $k = 4$ .

$S = 6.849$      $R\text{-Sq} = 97.1\%$      $R\text{-Sq(adjusted)} = 96.8\%$

**Coefficient of determination:**  
 $1 - \text{SSE}/\text{SST}$ , proportion of observed Y variation explained by the regression model

**Sequential Sums of squares:**  
 conditional contribution of each independent variable to SSR given the variables already entered into the model.

Source	DF	SS	MS	F	P
Regression	4	15913.0			
Residual Error	10	469.1			
Total	14	16382.2			

Source	DF	Seq SS
SqFeet	1	14829.3
NumFlrs	1	0.9
Bdrms	1	166.4
Baths	1	916.5

Slide 15

FOR6520, AUU, Ivo Dinov

## Testing the Usefulness of the Model

- The first question to ask is whether the regression model is of any use in predicting  $y$ .
- If it is not, then the value of  $y$  does not change, regardless of the value of the independent variables,  $x_1, x_2, \dots, x_k$ . This implies that the partial regression coefficients,  $\beta_1, \beta_2, \dots, \beta_k$  are all zero.

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0 \text{ versus } H_a : \text{at least one } \beta_i \text{ is not zero}$$

Slide 16

FOR6520, AUU, Ivo Dinov

## The F Test

- You can test the overall usefulness of the model using an F test. If the model is useful, MSR will be large compared to the unexplained variation, MSE.

To test  $H_0$  : model is useful in predicting  $y$  is equivalent to

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$\text{TestStatistic } F = \frac{\text{MSR}}{\text{MSE}}$$

Reject  $H_0$  if  $F > F_{\alpha}$  with  $k$  and  $n - k - 1$  df.

Slide 17

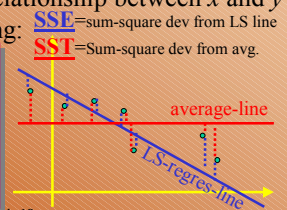
FOR6520, AUU, Ivo Dinov

## Measuring the Strength of the Relationship

- If the independent variables are useful in predicting  $y$ , you will want to know how well the model fits.
- The strength of the relationship between  $x$  and  $y$  can be measured using:

$R^2$  = Multiple coefficient of determination :

$$R^2 = 1 - \frac{\text{SSE}}{\text{SST}}$$



Slide 18

FOR6520, AUU, Ivo Dinov

## Measuring the Strength of the Relationship

- Since Total SS = SSR + SSE,  $R^2$  measures the proportion of the total variation in the responses that can be explained by using the independent variables in the model.
- ✓ the percent reduction the total variation by using the regression equation rather than just using the sample mean  $\bar{y}$  to estimate  $y$ .

$$R^2 = \frac{\text{SSR}}{\text{Total SS}} \quad \text{and} \quad F = \frac{\text{MSR}}{\text{MSE}} = \frac{R^2 / k}{(1 - R^2) / (n - k - 1)}$$

Slide 19

FOR6520, AUU, Joe Dimov

## Testing the Partial Regression Coefficients

- Is a particular independent variable useful in the model, *in the presence of all the other independent variables*? The test statistic is function of  $b_i$ , our best estimate of  $\beta_i$ .

$$H_0 : \beta_i = 0 \text{ versus } H_a : \beta_i \neq 0$$

$$\text{Test statistic : } t = \frac{b_i - 0}{\text{SE}(b_i)}$$

which has a  $t$  distribution with error  $df = n - k - 1$ .

Slide 20

FOR6520, AUU, Joe Dimov

## The Real Estate Problem

Is the overall model useful in predicting list price? How much of the overall variation in the response is explained by the regression model?

S = 6.849      R-Sq = 97.1%      R-Sq(adj) = 96.0%					
Analysis of Variance					
Source	DF	SS	MS	F	P
Regression	3	3978.3	1326.1	84.80	0.000
Error	9	45.9	5.1		
Total	12				

$R^2 = .971$  indicates that 97.1% of the overall variation is explained by the regression model.

$F = \text{MSR}/\text{MSE} = 84.80$  with  $p\text{-value} = .000$  is highly significant. The model is very useful in predicting the list price of homes.

Slide 21

FOR6520, AUU, Joe Dimov

## The Real Estate Problem

To test  $H_0 : \beta_3 = 0$ , the test statistic is  $t = -0.59$  with  $p\text{-value} = .565$ .

The  $p\text{-value}$  is larger than .05 and  $H_0$  is not rejected.

We cannot conclude that number of bedrooms is a valuable predictor in the presence of the other variables.

Perhaps the model could be refit without  $x_3$ .

Regress The List Price					
Predicted List Price = -2.67 Bdrms + 30.3 Baths					
				P	
Constant	10.703	7.207	2.04	0.069	
SqFeet	6.2698	0.7252	8.65	0.000	
NumFlrs	-16.203	6.212	-2.61	0.026	
Bdrms	-2.673	4.494	-0.59	0.565	
Baths	30.271	6.849	4.42	0.001	

Slide 22

FOR6520, AUU, Joe Dimov

## Comparing Regression Models

- The strength of a regression model is measured using  $R^2 = \text{SSR}/\text{Total SS}$ . This value will only increase as variables are added to the model.
- To fairly compare two models, it is better to use a measure that has been adjusted using  $df$ :

$$R^2(\text{adj}) = \left(1 - \frac{\text{MSE}}{\text{Total SS}/(n-1)}\right) 100\%$$

Slide 23

FOR6520, AUU, Joe Dimov

## Comparing Regression Models

- Remember that the results of a regression analysis are only valid when the necessary assumptions have been satisfied.

- ✓  $\epsilon$  are independent
- ✓ Have a mean 0 and common variance  $\sigma^2$  for any set  $x_1, x_2, \dots, x_k$ .
- ✓ Have a normal distribution.

Slide 24

FOR6520, AUU, Joe Dimov



## Diagnostic Tools

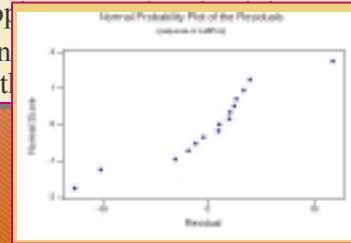
1. Normal probability plot of residuals
2. Plot of residuals versus fit or residuals versus variables

Slide 25

FOR6520, AU, Lec, Diagnostics

## Normal Probability Plot

- ✓ If the normality assumption is valid, the plot should resemble a straight line, slope is the standard deviation of the residuals.
- ✓ If not, the plot will show a non-linear pattern, indicating a failure of the normality assumption.

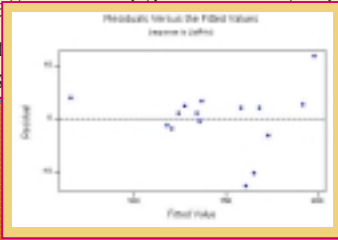


Slide 26

FOR6520, AU, Lec, Diagnostics

## Residuals versus Fits

- ✓ If the equal variance assumption is valid, the plot should appear as a random scatter of points around the zero line.
- ✓ If not, the plot will show a non-random pattern, indicating a failure of the equal variance assumption.



Slide 27

FOR6520, AU, Lec, Diagnostics

## Estimation and Prediction

- Once you have
  - ✓ determined that the regression line is useful
  - ✓ used the diagnostic plots to check for violation of the regression assumptions.
- You are ready to use the regression line to
  - ✓ Estimate the average value of  $y$  for a given value of  $x$
  - ✓ Predict a particular value of  $y$  for a given value of  $x$ .

Slide 28

FOR6520, AU, Lec, Diagnostics

## Estimation and Prediction

- Enter the appropriate values of  $x_1, x_2, \dots, x_k$  in Minitab. Minitab calculates

$$\hat{y} = b_0 + b_1x_1 + b_2x_2 + \dots + b_kx_k$$

- and both the confidence interval and the prediction interval.
- Particular values of  $y$  are more difficult to predict, requiring a wider range of values in the prediction interval.

Slide 29

FOR6520, AU, Lec, Diagnostics

## The Real Estate Problem

- Estimate the average list price for a home with 1000 square feet of living space and two baths with a 95% confidence interval.

We estimate that the average list price will be between \$110,860 and \$124,700 for a home like this.

Predicted Values for New Observations				
New Obs	Fit	SE Fit	95.0% CI	95.0% PI
1	117.78	3.11	( 110.86, 124.70)	( 101.02, 134.54)

Values of Predictors for New Observations				
New Obs	SqFeet	NumFlrs	Bdrms	Baths
1	10.0	1.00	3.00	2.00

Slide 30

FOR6520, AU, Lec, Diagnostics

## Using Regression Models

When you perform multiple regression analysis, use a step-by-step approach:

1. Obtain the fitted prediction model.
2. Use the analysis of variance  $F$  test and  $R^2$  to determine how well the model fits the data.
3. Check the  $t$  tests for the partial regression coefficients to see which ones are contributing significant information in the presence of the others.
4. If you choose to compare several different models, use  $R^2(\text{adj})$  to compare their effectiveness.
5. Use diagnostic plots to check for violation of the regression assumptions.

Slide 31

FOR6520, AUU, Jon Dimov

## A Polynomial Model

- A response  $y$  is related to a single independent variable  $x$ , but not in a linear manner. The polynomial model is:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_k x^k + \varepsilon$$

- When  $k = 2$ , the model is **quadratic**:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$$

- When  $k = 3$ , the model is **cubic**:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \varepsilon$$

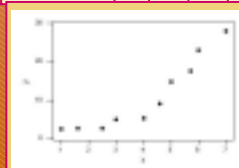
Slide 32

FOR6520, AUU, Jon Dimov

## Example

- A market research firm has observed the sales ( $y$ ) as a function of mass media advertising expenses ( $x$ ) for 10 different companies selling a similar product.

Company	1	2	3	4	5	6	7	8	9	10
Expenditure, $x$	1.0	1.6	2.5	3.0	4.0	4.6	5.0	5.7	6.0	7.0
Sales, $y$	2.5	2.6	2.7	5.0	5.3	9.1	14.8	17.5	23.0	28.0



Since there is only one independent variable, you could fit a linear, quadratic, or cubic polynomial model. Which would you pick?

Slide 33

FOR6520, AUU, Jon Dimov

## Two Possible Choices

A straight line model:  $y = \beta_0 + \beta_1 x + \varepsilon$

A quadratic model:  $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$

Overall  $F$  test is highly significant, as is the  $t$ -test of the slope.  $R^2 = .856$  suggests a good fit. Let's check the residual plots...

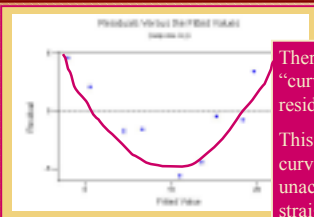
Here is the computer printout

Regression Analysis: y versus x					
The regression equation is					
$y = -6.47 + 4.34 x$					
Predictor	Coef	SE Coef	T	P	
Constant	-6.465	2.795	-2.31	0.049	
x	4.3355	0.6274	6.91	0.000	
S = 3.725	R-Sq = 85.6%		R-Sq(adj) = 83.9%		
Analysis of Variance					
Source	DF	SS	MS	F	P
Regression	1	662.46	662.46	47.74	0.000
Residual Error	8	111.00	13.88		
Total	9	773.46			

Slide 34

FOR6520, AUU, Jon Dimov

## Example



Fit the quadratic model:  
 $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$

There is a strong pattern of a "curve" leftover in the residual plot.

This indicates that there is a curvilinear relationship unaccounted for by your straight line model. You should have used the quadratic model!

Slide 35

FOR6520, AUU, Jon Dimov

## The Quadratic Model

Regression Analysis: y versus x, x-sq

The regression equation is

$$y = 4.66 - 3.03 x + 0.939 x\text{-sq}$$

Predictor	Coef	SE Coef	T	P
Constant	4.657	2.443	1.91	0.098
x	-3.030	1.395	-2.17	0.067
x-sq	0.9389	0.1739	5.40	0.001
S = 1.752	R-Sq = 97.2%		R-Sq(adj) = 96.4%	

Overall  $F$  test is highly significant, as is the  $t$ -test of the quadratic term  $\beta_2$ .  $R^2 = .972$  suggests a very good fit.

Let's compare the two models, and check the residual plots.

MS	F	P
375.99	122.49	0.000
3.07		

Slide 36

FOR6520, AUU, Jon Dimov

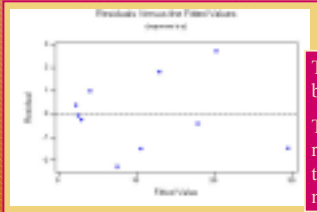
## Which Model to Use?

Use  $R^2(\text{adj})$  to compare the models:

The straight line model:  $y = \beta_0 + \beta_1 x + \varepsilon$   $R^2(\text{adj}) = 83.9\%$

The quadratic model:  $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$

$R^2(\text{adj}) = 96.4\%$



The quadratic model is better.

There are no patterns in the residual plot, indicating that this is the correct model for the data.

Slide 37

FOR6520, AU, Lec 10, Dimes

## Using Qualitative Variables

- Multiple regression requires that the response  $y$  be a quantitative variable.
- Independent variables can be either quantitative or qualitative.
- Qualitative variables** involving  $k$  categories are entered into the model by using  $k-1$  **dummy variables**.
- Example:** To enter **gender** as a variable, use  $x_i = 1$  if male; 0 if female

Slide 38

FOR6520, AU, Lec 10, Dimes

## Example

- Data was collected on 6 male and 6 female assistant professors. The researchers recorded their salaries ( $y$ ) along with years of experience ( $x_1$ ). The professor's gender enters into the model as a dummy variable:  $x_2 = 1$  if male; 0 if not.

Professor	Salary, $y$	Experience, $x_1$	Gender, $x_2$	Interaction, $x_1 x_2$
1	\$50,710	1	1	1
2	49,510	1	0	0
...	...	...	...	...
11	55,590	5	1	5
12	53,200	5	0	0

Slide 39

FOR6520, AU, Lec 10, Dimes

## Example

- We want to predict a professor's salary based on years of experience and gender. We think that there may be a difference in salary depending on whether you are male or female. The model we choose includes experience ( $x_1$ ), gender ( $x_2$ ), and an interaction term ( $x_1 x_2$ ) to allow salary's for males and females to behave differently.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \varepsilon$$

Slide 40

FOR6520, AU, Lec 10, Dimes

## Output

What is the regression

Is the overall model useful in predicting  $y$ ?

**Regression Analysis: y versus x1, x2, x1x2**  
The regression equation is:  
 $y = 48593 + 969 x_1 + 867 x_2 + 260 x_1 x_2$

Predictor	Coef	SE Coef
Constant	48593.0	207.9
x1	969.00	63.67
x2	866.7	305.3
x1x2	260.13	87.06

The overall  $F$  test is  $F = 346.24$  with  $p$ -value = .000. The value of  $R^2 = .992$  indicates that the model fits very well.

Is there a difference in the relationship between salary and years of experience, depending on the gender of the professor?

Yes. The individual  $t$ -test for the interaction term is  $t = 2.99$  with  $p$ -value = .017. This indicates a significant interaction between gender and years of experience.

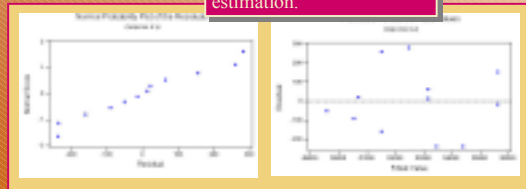
Slide 41

FOR6520, AU, Lec 10, Dimes

Have any of the re  
violated, or have v

It does not appear from the diagnostic plots that there are any violations of assumptions.

The model is ready to be used for prediction or estimation.



Slide 42

FOR6520, AU, Lec 10, Dimes

### Testing Sets of Parameters

- Suppose the demand  $y$  may be related to five independent variables, but that the cost of measuring three of them is very high.
- If it could be shown that these three contribute little or no information, they can be eliminated.
- You want to test the null hypothesis
- $H_0 : \beta_3 = \beta_4 = \beta_5 = 0$ —that is, the independent variables  $x_3$ ,  $x_4$ , and  $x_5$  contribute no information for the prediction of  $y$ —versus the alternative hypothesis:
- $H_a$  : At least one of the parameters  $\beta_3$ ,  $\beta_4$ , or  $\beta_5$  differs from 0—that is, at least one of the variables  $x_3$ ,  $x_4$ , or  $x_5$  contributes information for the prediction of  $y$ .

Slide 43

FOR6520, AUU, Ivo Dinov

### Testing Sets of Parameters

- To explain how to test a hypothesis concerning a set of model parameters, we define two models:
- Model One (reduced model)**  

$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_r x_r$$
- Model Two (complete model)**  

$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_r x_r + \beta_{r+1} x_{r+1} + \beta_{r+2} x_{r+2} + \cdots + \beta_k x_k$$
  - terms in model 1      additional terms in model 2

Slide 44

FOR6520, AUU, Ivo Dinov

### Testing Sets of Parameters

- The test of the hypothesis
- $H_0 : \beta_3 = \beta_4 = \beta_5 = 0$
- $H_a$  : At least one of the  $\beta_i$  differs from 0
- uses the test statistic  $F = \frac{(SSE_1 - SSE_2)/(k - r)}{MSE_2}$

where  $F$  is based on  $df_1 = (k - r)$  and  $df_2 = n - (k + 1)$ .

The rejection region for the test is identical to other analysis of variance  $F$  tests, namely  $F > F_{\alpha}$ .

Slide 45

FOR6520, AUU, Ivo Dinov

### Stepwise Regression

- A stepwise regression analysis fits a variety of models to the data, adding and deleting variables as their significance in the presence of the other variables is either **significant** or **nonsignificant**, respectively.
- Once the program has performed a sufficient number of iterations and no more variables are significant when added to the model, and none of the variables are nonsignificant when removed, the procedure stops.
- These programs **always fit first-order models** and are not helpful in detecting curvature or interaction in the data.

Slide 46

FOR6520, AUU, Ivo Dinov

### Some Points of Caution

- ✓ **Causality:** Be careful not to deduce a causal relationship between a response  $y$  and a variable  $x$ .
- ✓ **Multicollinearity:** Neither the size of a regression coefficient nor its  $t$ -value indicates the importance of the variable as a contributor of information. This may be because two or more of the predictor variables are highly correlated with one another; this is called **multicollinearity**.

Slide 47

FOR6520, AUU, Ivo Dinov

### Key Concepts

#### I. The General Linear Model

- $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k + \varepsilon$
- The random error  $\varepsilon$  has a normal distribution with mean 0 and variance  $\sigma^2$ .

#### II. Method of Least Squares

- Estimates  $b_0, b_1, \dots, b_k$  for  $\beta_0, \beta_1, \dots, \beta_k$  are chosen to minimize SSE, the sum of squared deviations about the regression line  $\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \cdots + b_k x_k$ .
- Least-squares estimates are produced by computer.

Slide 50

FOR6520, AUU, Ivo Dinov



## Key Concepts

### III. Analysis of Variance

1. Total SS = SSR + SSE, where Total SS =  $S_{yy}$ .  
The ANOVA table is produced by computer.
2. Best estimate of  $\sigma^2$  is

$$MSE = \frac{SSE}{n - k - 1}$$

### IV. Testing, Estimation, and Prediction

1. A test for the significance of the regression,  
 $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$ , can be implemented using the analysis of variance  $F$  test:

$$F = \frac{MSR}{MSE}$$

Slide 51

FOR6520, AUU, Ivo Dinov

## Key Concepts

2. The strength of the relationship between  $x$  and  $y$  can be measured using

$$R^2 = 1 - \frac{SST}{SSE}$$

which gets closer to 1 as the relationship gets stronger.

3. Use residual plots to check for non-normality, inequality of variances, and an incorrectly fit model.
4. Significance tests for the partial regression coefficients can be performed using the Student's  $t$  test with error  $df = n - k - 1$ :

$$t = \frac{b_i - \beta_i}{SE(b_i)}$$

Slide 52

FOR6520, AUU, Ivo Dinov

## Key Concepts

5. Confidence intervals can be generated by computer to estimate the average value of  $y$ ,  $E(y)$ , for given values of  $x_1, x_2, \dots, x_k$ . Computer-generated prediction intervals can be used to predict a particular observation  $y$  for given value of  $x_1, x_2, \dots, x_k$ . For given  $x_1, x_2, \dots, x_k$ , prediction intervals are always wider than confidence intervals.

Slide 53

FOR6520, AUU, Ivo Dinov

## Key Concepts

### V. Model Building

1. The number of terms in a regression model cannot exceed the number of observations in the data set and should be considerably less!
2. To account for a curvilinear effect in a **quantitative** variable, use a second-order polynomial model. For a cubic effect, use a third-order polynomial model.
3. To add a **qualitative** variable with  $k$  categories, use  $(k - 1)$  dummy or indicator variables.
4. There may be interactions between two qualitative variables or between a quantitative and a qualitative variable. Interaction terms are entered as  $\beta_{xy}x_j$ .
5. Compare models using  $R^2(\text{adj})$ .

Slide 54

FOR6520, AUU, Ivo Dinov