



Management Science

Regression

Dan Brandon, Ph.D., PMP

Session Objectives

- Identify variables and use them in a regression model
- Develop simple linear regression equations from sample data and interpret the slope and intercept
- Compute the coefficient of determination and the coefficient of correlation and interpret their meanings
- List the assumptions used in regression and use residual plots to identify problems
- Develop a multiple regression model and use it to predict
- Use dummy variables to model categorical data
- Determine which variables should be included in a multiple regression model
- Transform a nonlinear function into a linear one
- Understand and avoid common mistakes made in the use of regression analysis

Goals Of Regression Analysis

- To determine if two or more variables are **related** and to what degree
- In situations where one variable precedes another in time, we maybe able to use regression analysis to **predict** the value of one variable by knowing the value of the other variable
- **What might be some business variables that are related ?**

Regression and Decision Making

- Business (and personal) decisions are often based on the perceived relationship between two variables
- For example, a manager might decide how much to spend on advertising based on his believe about the relationship between advertising and sales revenue
- You might decide how much time to spend studying based on your believe about the relationship between study time and test grade
- You might spend the time and money to get a college degree based on your believe about the relationship between education and lifetime income

Regression Analysis Concepts

- Regression analysis can be used to quantify how variables are **related** or **associated** (or “**covary**” or are “**correlated**”)
- It may be the case that one variable’s (x) occurrence precedes the other variable’s occurrence (y) in time
- Regression analysis may allow us to **predict** the value of one variable based on the observed value of another
- But we must be careful to note that ‘**correlation does not imply causality**’
- Therefore, regression analysis *alone* cannot be used to conclude causality – that is, we cannot conclude from regression analysis *alone* that one variable’s outcome **caused** another variable’s outcome !

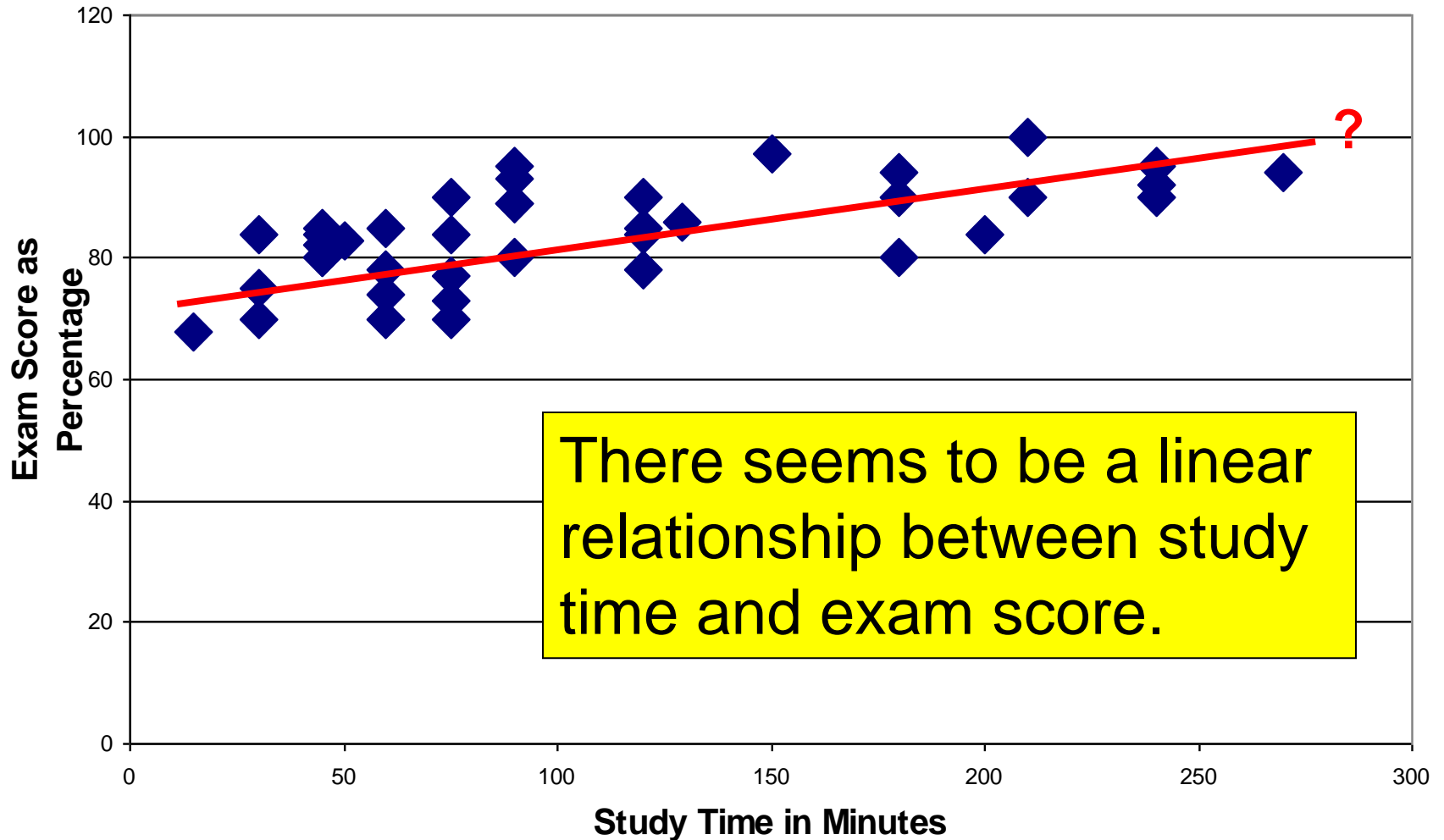
Relationship Between Variables

- Let's say that a professor at a college believes there is a *relationship* between the time that a student studies for an exam and that student's exam score
- The first step that should be made is to develop a **scatter plot** - plotting the values of students' reported study times and their exam scores to see if there appears to be a linear* relationship

* We may ultimately discover that the relationship is not linear but *curvilinear*



Study Time vs. Exam Score



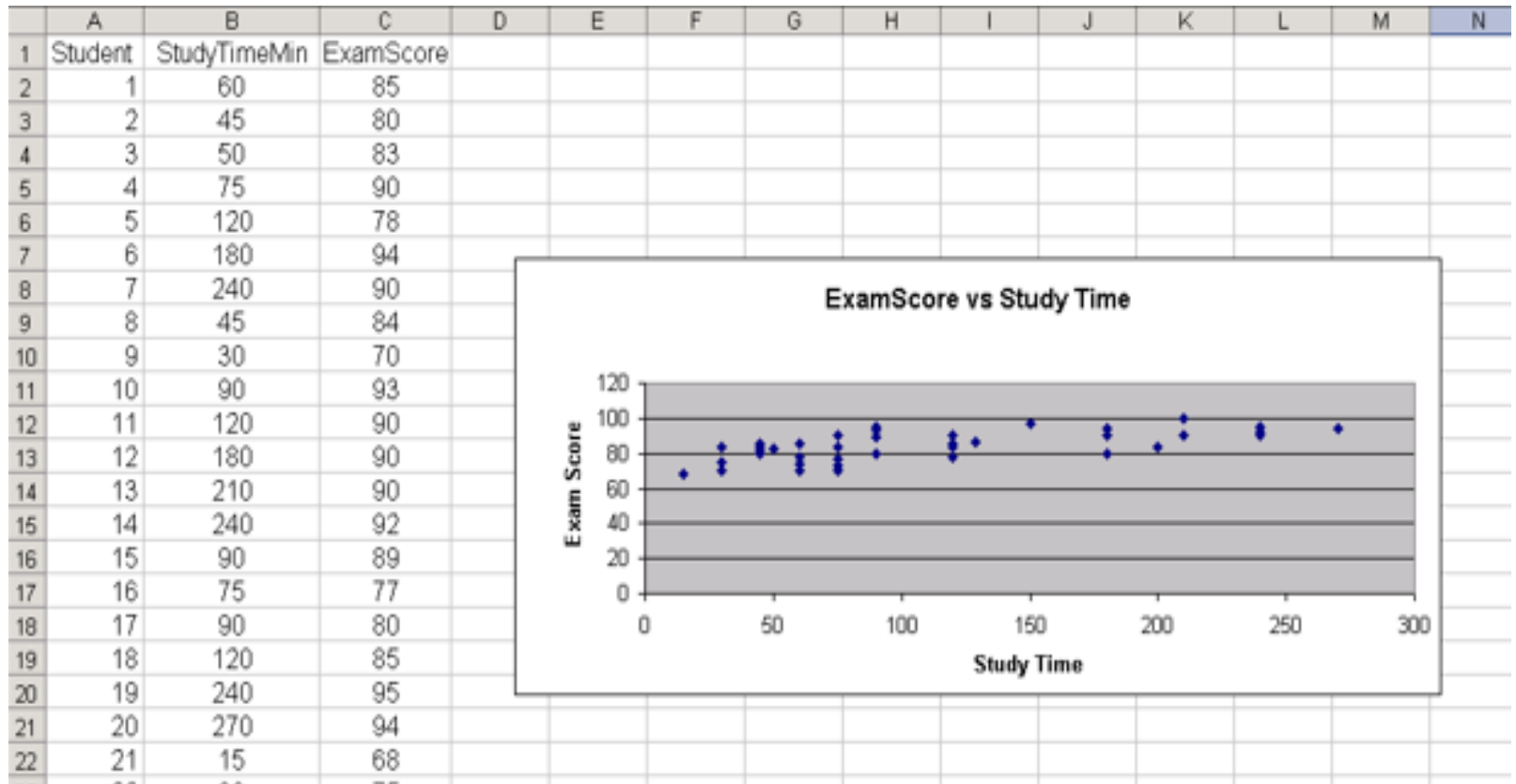
Scatter Plots in Excel



	A	B	C	D	E	F	G	H	I	J
1	Student	StudyTimeMin	ExamScore							
2	1	60	85							
3	2	45	80							
4	3	50	83							
5	4	75	90							
6	5	120	78							
7	6	180	94							
8	7	240	90							
9	8	45	84							
10	9	30	70							
11	10	90	93							
12	11	120	90							
13	12	180	90							
14	13	210	90							
15	14	240	92							
16	15	90	89							
17	16	75	77							
18	17	90	80							
19	18	120	85							
20	19	240	95							
21	20	270	94							
22	21	15	68							
23	22	30	75							
24	23	45	85							
25	24	90	89							
26	25	75	73							

Completed Excel Scatter Chart

[insert -> charts -> scatter (x,y)]



The Linear Correlation Coefficient

- Although the scatter plot indicates that there is somewhat of a linear relationship between study time and exam score, we can use a more precise measure to quantify the degree of association
- The **linear correlation coefficient (r)** measures the strength of the linear relationship between the metric values in the sample
- The linear correlation coefficient is also referred to as the **Pearson product moment correlation coefficient** after Karl Pearson, the person who developed this measure - Some just call it "***Pearson's r*** "



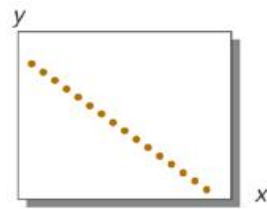
Pearson's r



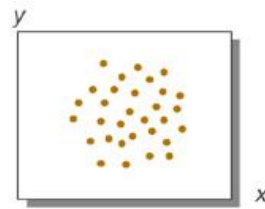
- **Pearson's r** measures the strength of the association between two variables
- Its value can be any number from -1 (indicating a perfect, negative correlation) to +1 (indicating a perfect, positive correlation)
 - See graphical examples on a following slide →



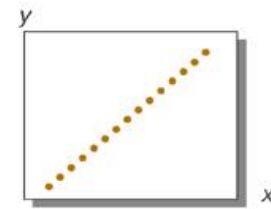
Scatterplots of Relationships and r Values



$r = -1.00$



$r = 0$



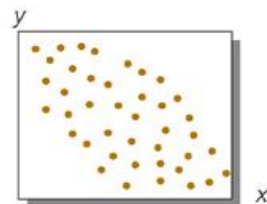
$r = +1.00$



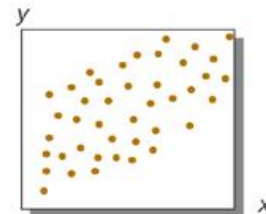
$r = -.90$



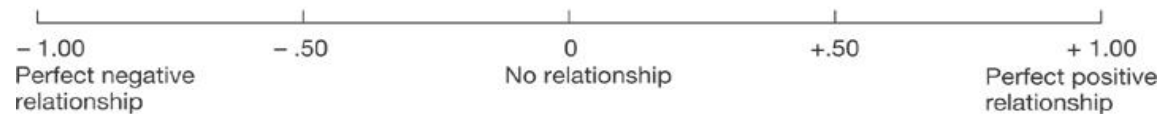
$r = +.90$



$r = -.40$



$r = +.40$



An example:

	Student	StudyTime Min	ExamScore	
1	1	60	85	
2	2	45	80	
3	3	50	83	
4	4	75	90	
5	5	120	78	
6	6	180	94	
7	7	240	90	
8	8	45	84	
9	9	30	70	
10	10	90	93	
11	11	120	90	
12	12	180	90	
13	13	210	90	
14	14	240	92	
15	15	90	89	
16	16	75	77	
17	17	90	80	
18	18	120	85	
19	19	240	95	
20	20	270	94	
21	21	15	68	
22	22	30	75	
23	23	45	85	
24	24	90	89	
25	25	75	73	
26	26	60	78	
27	27	60	74	
28	28	120	84	
29	29	45	80	
30	30	20	84	

Data View Variable View

start

Untitled - SP

Let's say this is the data set that contains 40 students' reported study times and exam scores

Let's refer to study time as 'x' and exam score as 'y'



The Calculation of r (Part 1)

First, we calculate $\Sigma(x)(y)$, Σx^2 , and Σy^2

Measuring the Association
between study time and exam score

StudyTime in Minutes (x)	Exam Score as Percentage (y)	(x) * (y)	x^2	y^2
60	85	5,100	3,600	7,225
45	80	3,600	2,025	6,400
200	84	16,800	40,000	7,056
210	100	21,000	44,100	10,000
4,354	3,367	380,809	666,416	286,001

Skipped
rows

The Calculation of r (con't)

Next, we calculate r - the (standardized) correlation coefficient

$$r = \frac{n (\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{n (\Sigma x^2) - (\Sigma x)^2} * \sqrt{n(\Sigma y^2) - (\Sigma y)^2}}$$

$$r = \frac{40 (380,809) - (4,354)(3,367)}{\sqrt{40 (666,416) - (4,354)^2} * \sqrt{n(286,001) - (3,367)^2}}$$

r = + .642 r measures the strength of the association

Interpreting r

- The general rule is that the absolute value of r should be $> .7$ if it is to be considered significant
- But the specific 'significance' criteria is related to sample size
 - The smaller the sample size, the larger r must be to be considered significant

Pearson's r (correlation) in Excel

Microsoft Excel - DataSetsRM.xls

File Edit View Insert Format Tools Data Window Help

Type a question for help

STDEV

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	Student	StudyTimeMin	ExamScore											
2	1	60	85											
3	2	45	80											
4	3	50	83											
5	4	75	90											
6	5	120	78											
7	6	180	94			Pearson's r:	=							
8	7	240	90											
9	8	45	84											
10	9	30	70											
11	10	90	93											
12	11	120	90											
13	12	180	90											
14	13	210	90											
15	14	240	92											
16	15	90	89											
17	16	75	77											
18	17	90	80											
19	18	120	85											
20	19	240	95											
21	20	270	94											
22	21	15	68											
23	22	30	75											
24	23	45	85											
25	24	90	89											
26	25	75	73											

Insert Function

Search for a function:

Type a brief description of what you want to do and then click Go

Go

Or select a category: Statistical

Select a function:

- BINOMDIST
- CHIDIST
- CHIINV
- CHITEST
- CONFIDENCE
- CORREL
- COUNT

CORREL(array1,array2)

Returns the correlation coefficient between two data sets.

[Help on this function](#)

OK Cancel

Auto Ratings Workers Website Usability QualityExam Salaries pizza StudyTime Butle

Draw AutoShapes

Edit

NUM

Start

C:\Courses\... Chapter 19... 1695_19b.ppt Microsoft E... Using Excel t... Norton

12:06 PM



Using correlation function:

The screenshot shows a Microsoft Excel spreadsheet with the following data:

Student	StudyTimeMin	ExamScore
1	60	85
2	45	80
3	50	83
4	75	90
5	120	78
6	180	94
7	240	90
8	45	84
9	30	70
10	90	93
11	120	90
12	180	90
13	210	90
14	240	92
15	90	89
16	75	77
17	90	80
18	120	85
19	240	95
20	270	94
21	15	68
22	30	75
23	45	85
24	90	89
25	75	73

The formula bar shows the formula: `=CORREL(B2:B41,C2:C41)`. The function arguments dialog box is open, showing the following details:

- Function: CORREL
- Array1: B2:B41 (values: {60;45;50;75;120;180;240;45;30;90;120;180;210;240;90;75;90;120;240;270;15;30;45;90;75})
- Array2: C2:C41 (values: {85;80;83;90;78;94;90;84;70;93;90;90;90;92;89;77;80;85;95;94;68;75;85;89;73})
- Formula result: 0.641723363

The dialog box also includes a description: "Returns the correlation coefficient between two data sets." and a note: "Array2 is a second cell range of values. The values should be numbers, names, arrays, or references that contain numbers."

Result of function calculation:

Microsoft Excel - DataSetsRM.xls

File Edit View Insert Format Tools Data Window Help

Type a question for help

Reply with Changes... End Review...

Security...

F7 =CORREL(B2:B41,C2:C41)

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	Student	StudyTimeMin	ExamScore											
2	1	60	85											
3	2	45	80											
4	3	50	83											
5	4	75	90											
6	5	120	78											
7	6	180	94		Pearson's r:	0.6417								
8	7	240	90											
9	8	45	84											
10	9	30	70											
11	10	90	93											
12	11	120	90											
13	12	180	90											
14	13	210	90											
15	14	240	92											
16	15	90	89											
17	16	75	77											
18	17	90	80											
19	18	120	85											
20	19	240	95											
21	20	270	94											
22	21	15	68											
23	22	30	75											
24	23	45	85											
25	24	90	89											
26	25	75	73											

Auto Ratings Workers WebSite Usability QualityExam Salaries pizza StudyTime Butle

Draw AutoShapes

Ready NUM

Start C:\Courses\... Chapter 19... 695_19b.ppt Microsoft E... Norton 12:09 PM

Or Using the Pearson function:

The screenshot shows a Microsoft Excel window titled "DataSetsRM.xls". The spreadsheet contains a table with the following data:

Student	StudyTimeMin	ExamScore
1	60	85
2	45	80
3	50	83
4	75	90
5	120	78
6	180	94
7	240	90
8	45	84
9	30	70
10	90	93
11	120	90
12	180	90
13	210	90
14	240	92
15	90	89
16	75	77
17	90	80
18	120	85
19	240	95
20	270	94
21	15	68
22	30	75
23	45	85
24	90	89
25	75	73

In cell F7, the text "Pearson's r:" is followed by an equals sign and a blank cell. The "Insert Function" dialog box is open, showing the "Statistical" category selected. The "PEARSON" function is highlighted in the list. The description of the function is: "Returns the Pearson product moment correlation coefficient, r."

The taskbar at the bottom shows the Start button, several open applications (Auto Ratings, Workers, WebSite Usability, QualityExam, Salaries, pizza, StudyTime, Butle), and the system clock showing 12:18 PM.

Pizza Sales – What is the R value ?

	A	B	C
1	PizzaStore	CampusSize	Sales
2	1	2	58
3	2	6	105
4	3	8	88
5	4	8	118
6	5	12	117
7	6	16	137
8	7	20	157
9	8	20	169
10	9	22	149
11	10	26	202

Independent variable ?

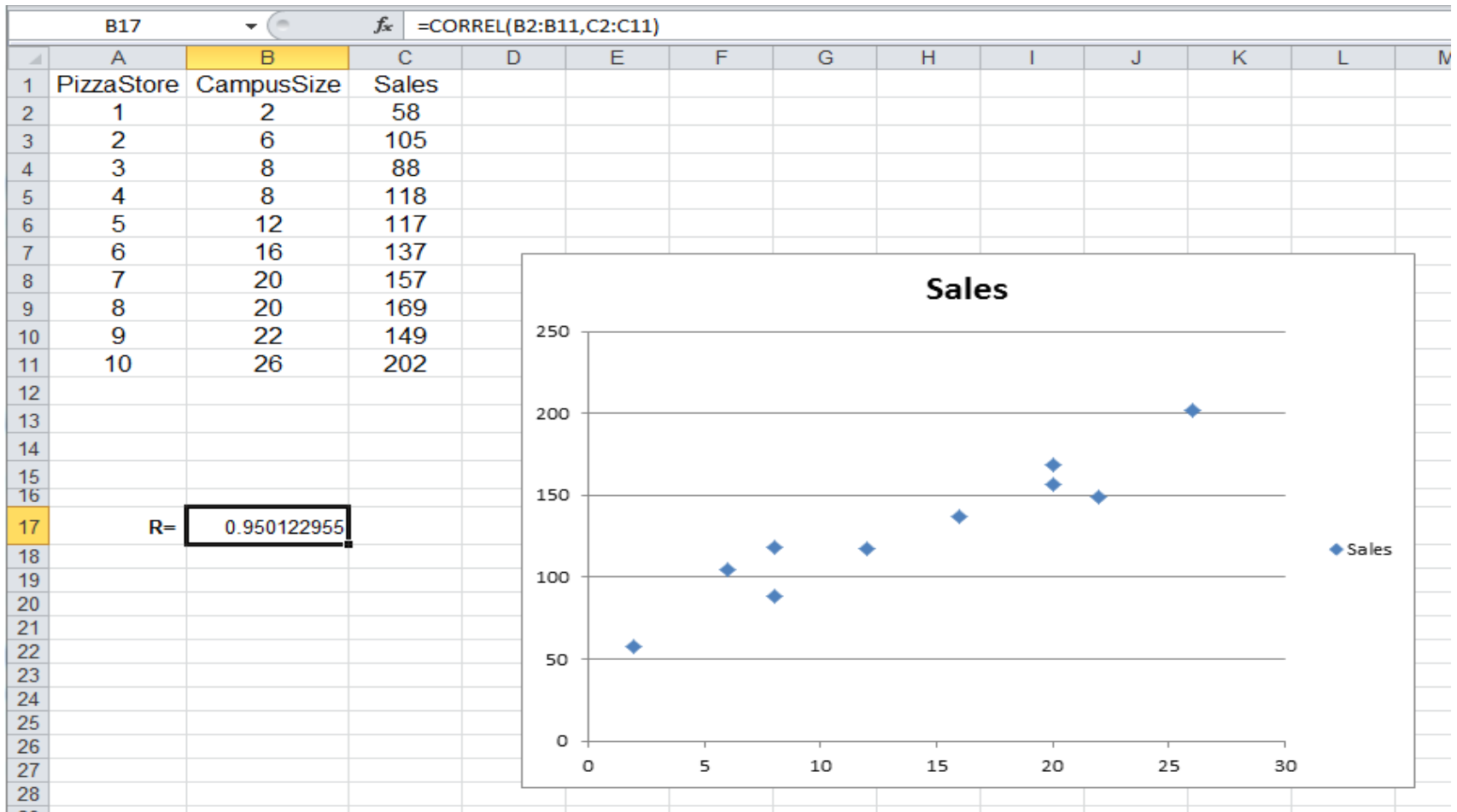
Draw a scatter plot also...

Wait....



Don't look ahead, until
you have your answer !

Pizza Sales (con't)

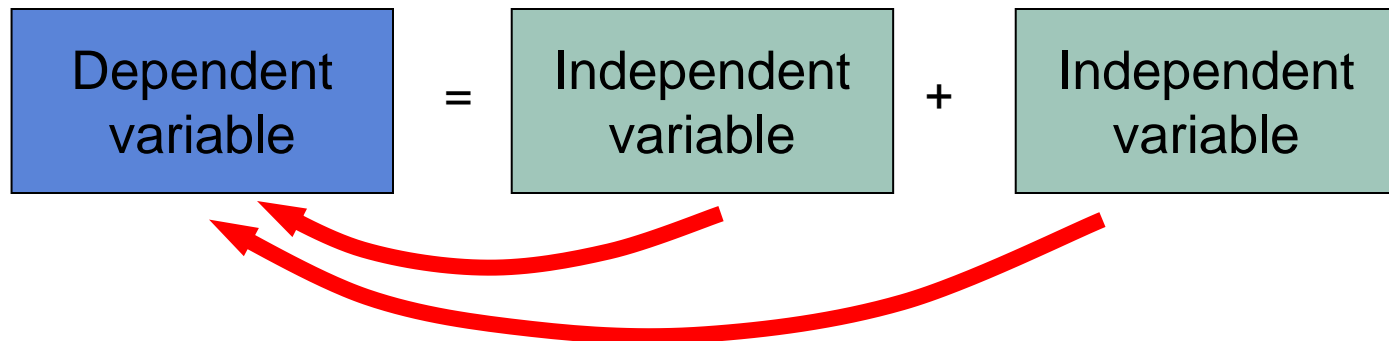


Simple Linear Regression

- Simple Linear Regression (SLR) allows us to explain and **possibly predict** the effect of one variable upon another
- For example one might want to analyze the effect of local unemployment rate on the local crime rate
- In this example the local crime rate is called the “dependent” variable (it depends on the other), and the unemployment rate is called the independent variable
- **“Simple” means there is only one independent variable that will be considered**

Independent & Dependent Variables

- The variable to be predicted is called the *dependent variable*
 - Sometimes called the *response variable*
- The value of this variable depends on the value of the *independent variable*
 - Sometimes called the *explanatory* or *predictor variable*



Sales and Advertising Data for Appleglo



Appleglo	First-Year Advertising Expenditures (\$ millions)	First-Year Sales (\$ millions)
Region	x	y
Maine	1.8	104
New Hampshire	1.2	68
Vermont	0.4	39
Massachusetts	0.5	43
Connecticut	2.5	127
Rhode Island	2.5	134
New York	1.5	87
New Jersey	1.2	77
Pennsylvania	1.6	102
Delaware	1.0	65
Maryland	1.5	101
West Virginia	0.7	46
Virginia	1.0	52
Ohio	0.8	33

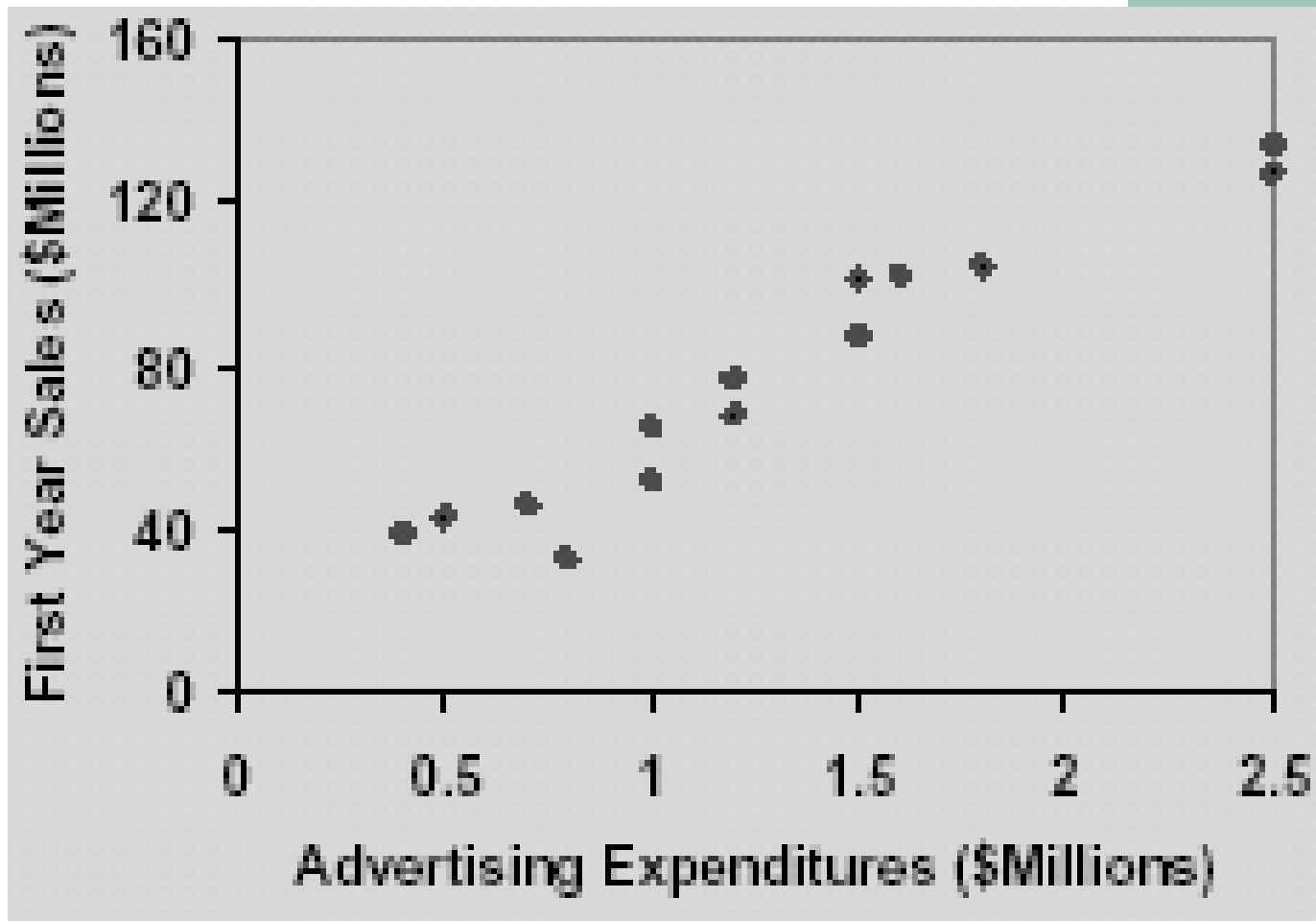
What is the independent variable, the dependent variable?

SLR Questions

- Business questions:
 - What is the relationship between advertising expenditure and sales ?
 - How much can we increase sales by a certain increase in advertising budget ?
 - How confident are we in our above estimate ?
- Here the dependent variable is sales and the independent variable is the advertising budget (sales depends upon advertising)



“Scatter Plot” of Sales versus Advertising



Simple Linear Regression

- True values for the slope and intercept are not known so they are estimated using sample data

$$\hat{Y} = b_0 + b_1X$$

where

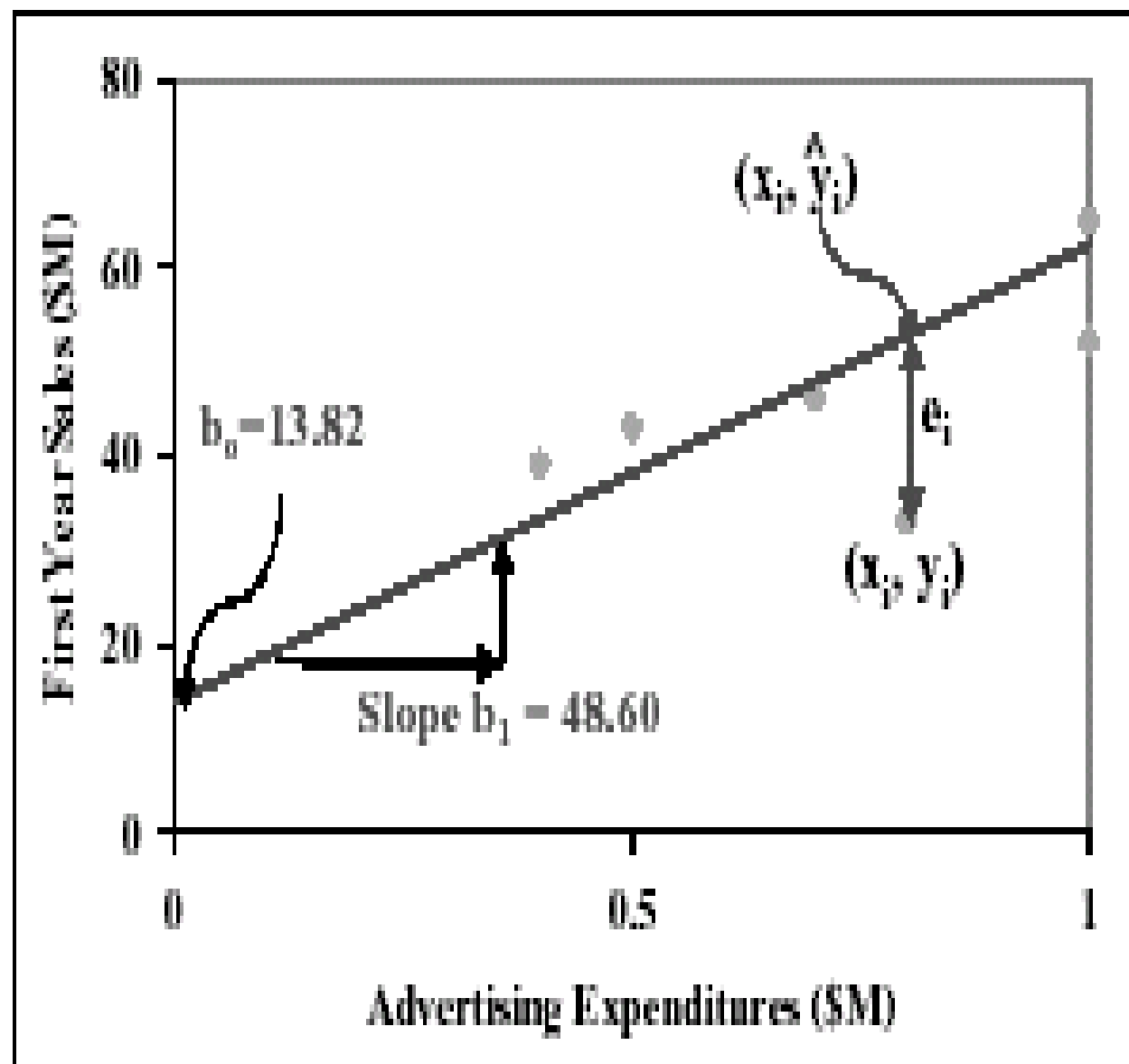
\hat{Y} = dependent variable (response)

X = independent variable (predictor or explanatory)

b_0 = intercept (value of Y when $X = 0$)

b_1 = slope of the regression line

How do we choose the line that “best” fits the data?



Best choices:

$$b_0 = 13.82$$

$$b_1 = 48.60$$

Prediction Formula

- The “prediction formula” here is:
 - $\text{Sales} = 13.82 + 48.6 * \text{Advertising}$
- Thus a dollar of advertising adds \$48.60 of sales
- We would have 13.82 million in sales even if we did not advertise
- The best fit regression line is the one that minimizes the sum of the distances of the data points from that line
- Or **minimizes the sum of the squares of the difference between the dependent variable sample values (y_i) and the predicted value** using the equation of the line
 - This formula is called the SS_{error} or “SSE”

Example: Study Time



Recall our data on the
'studyTime' worksheet



	A	B	C	D	E
1	Student	StudyTimeMin	ExamScore		
2	1	60	85		
3	2	45	80		
4	3	50	83		
5	4	75	90		
6	5	120	78		
7	6	180	94		
8	7	240	90		
9	8	45	84		
10	9	30	70		
11	10	90	93		
12	11	120	90		
13	12	180	90		
14	13	210	90		
15	14	240	92		
16	15	90	89		
17	16	75	77		
18	17	90	80		
19	18	120	85		
20	19	240	95		
21	20	270	94		
22	21	15	68		
23	22	30	75		
24	23	45	85		
25	24	90	89		
26	25	75	73		
27	26	60	78		
28	27	60	74		
29	28	120	84		
30	29	45	80		
31	30	30	84		
32	31	75	70		
33	32	90	95		
34	33	150	97		
35	34	129	86		
36	35	75	84		
37	36	60	70		
38	37	45	82		
39	38	180	80		
40	39	200	84		
41	40	210	100		
42					

Coefficients of the Regression Line

- The regression **formulas for the coefficients** are shown below, and we use the sample data to derive estimates of β_0 and β_1 (see appendix for full derivation):

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

Calculating the coefficients of the regression line

Microsoft Excel - studytime

File Edit View Insert Format Tools Data Window Help

Arial 12 B I U

B55 fx

	A	B	C	D	E	F	G	H
1	Using Regression to explain and predict exam score from study time							
2	Simple Linear Regression							
3		Study Time	Exam Score					
4	Student i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$	
5	1	60	85					
6	2	45	80					
7	3	50	83					
8	4	75	90					
9	5	120	78					
12	38	180	80					
13	39	200	84					
14	40	210	100					
15	Averages:			Totals:				
17	$\beta_1 = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2} =$							
19	$\beta_0 = \bar{y} - \beta_1 \bar{x} =$							
21	The regression equation ($\hat{y} = \beta_0 + \beta_1 x$):							

The coefficients can be calculated 'manually' in Excel

Follow along these next slides to see



G

Calculating the coefficients of the regression line

Microsoft Excel - studytime

File Edit View Insert Format Tools Data Window Help

Arial 12 B I U

TTEST X ✓ ✗ =AVERAGE(B5:B44)

	A	B	C	D			
1	Using Regression to explain and predict exam						
2	Simple Linear Regression						
3		Study Time	Exam Score				
4	Student i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
5	1	60	85	-48.85			
6	2	45	80				
7	3	50	83				
8	4	75	90				
9	5	120	78				
42	38	180	80				
43	39	200	84				
44	40	210	100				
45	=AVERAGE(B5:B44)			Totals:			
47	$\beta_1 = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2} =$						
49	$\beta_0 = \bar{y} - b_1\bar{x} =$						
51	The regression equation ($\hat{y} = \beta_0 + \beta_1 x$):						

Calculate (\bar{x}) , the overall average study time:

B45: = average(B5:B44)

Rows hidden

Calculating the coefficients of the regression line

Microsoft Excel - studytime

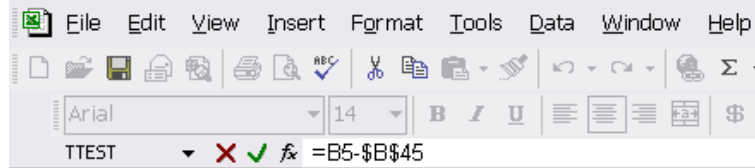
Using Regression to explain and predict exam						
Simple Linear Regression						
	Study Time	Exam Score				
Student i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	60	85	-48.85			
2	45	80				
3	50	83				
4	75	90				
5	120	78				
38	180	80				
39	200	84				
40	210	100				
Averages:	1	=AVERAGE(C5:C44)				
Totals:						
$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$						
$\beta_0 = \bar{y} - b_1\bar{x} =$						
The regression equation ($\hat{y} = \beta_0 + \beta_1 x$):						

Calculate (\bar{y}) , the overall average test score:

C45: = average(C5:C44)

Calculating the coefficients of the regression line

Microsoft Excel - studytime



Calculate $(x_i - \bar{x})$ for the first observation in the data set:

D5: =B5 - \$B\$45

Using Regression to explain and predict						
Simple Linear Regression						
	Study Time	Exam Score				
Student i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	60		=B5-\$B\$45			
2	45	80				
3	50	83				
4	75	90				
5	120	78				
38	180	80				
39	200	84				
40	210	100				
Averages:	108.85	84.18			Totals:	
$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$						
$\beta_0 = \bar{y} - b_1\bar{x} =$						
The regression equation ($\hat{y} = \beta_0 + \beta_1 x$):						

Calculating the coefficients of the regression line

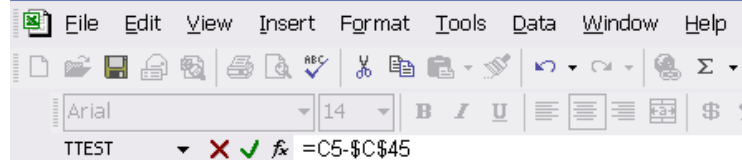
Microsoft Excel - studytime

Copy the formula in D5 down to D44 to calculate $(x_i - \bar{x})$ for the remaining observations in the data set

Using Regression to explain the relationship between Study Time and Exam Score						
Student i	Study Time x_i	Exam Score y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	60	85	-48.85			
2	45	80	-63.85			
3	50	83	-58.85			
4	75	90	-33.85			
5	120	78	11.15			
38	180	80	71.15			
39	200	84	91.15			
40	210	100	101.15			
Averages:	108.85	84.18				
$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$						
$\beta_0 = \bar{y} - b_1\bar{x} =$						
The regression equation ($\hat{y} = \beta_0 + \beta_1 x$):						

Calculating the coefficients of the regression line

Microsoft Excel - studytime



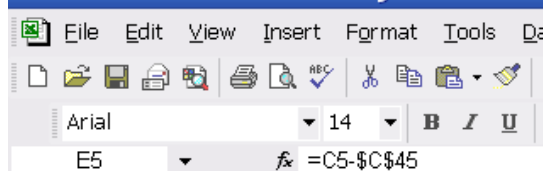
Calculate $(y_i - \bar{y})$ for the first observation in the data set:

E5: $=C5 - \$C\45

Using Regression to explain and predict						
Simple Linear Regression						
	Study Time	Exam Score				
Student i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	60	85		$=C5 - \$C\45		
2	45	80	-63.85			
3	50	83	-58.85			
4	75	90	-33.85			
5	120	78	11.15			
38	180	80	71.15			
39	200	84	91.15			
40	210	100	101.15			
Averages:	108.85	84.18		Totals:		
$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$						
$\beta_0 = \bar{y} - b_1\bar{x} =$						
The regression equation ($\hat{y} = \beta_0 + \beta_1 * x$):						

Calculating the coefficients of the regression line

Microsoft Excel - studytime



Copy the formula in E5 down to E44 to calculate $(y_i - \bar{y})$ for the remaining observations in the data set

Using Regression to explain and predict exam score from study time						
Simple Linear Regression						
	Study Time	Exam Score				
Student i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	60	85	-48.85	0.83		
2	45	80	-63.85	-4.18		
3	50	83	-58.85	-1.18		
4	75	90	-33.85	5.83		
5	120	78	11.15	-6.18		
38	180	80	71.15	-4.18		
39	200	84	91.15	-0.17		
40	210	100	101.15	15.83		
Averages:	108.85	84.18		Totals.		
$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$						
$\beta_0 = \bar{y} - b_1\bar{x} =$						
The regression equation ($\hat{y} = \beta_0 + \beta_1^* x$):						

Calculating the coefficients of the regression line

Microsoft Excel - studytime

Using Regression to explain and						
Simple Linear Regression						
	Student i	Study Time x_i	Exam Score y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$
1	1	60	85	-48.85	0.83	=D5*E5
2	2	45	80	-63.85	-4.18	
3	3	50	83	-58.85	-1.18	
4	4	75	90	-33.85	5.83	
5	5	120	78	11.15	-6.18	
6	38	180	80	71.15	-4.18	
7	39	200	84	91.15	-0.17	
8	40	210	100	101.15	15.83	
9	Averages:	108.85	84.18		Totals:	
10	$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$					
11	$\beta_0 = \bar{y} - b_1\bar{x} =$					
12	The regression equation ($\hat{y} = \beta_0 + \beta_1 * x$):					

Calculate $(x_i - \bar{x})(y_i - \bar{y})$ for the first observation in the data set:

$$F5: =D5 * E5$$

Calculating the coefficients of the regression line

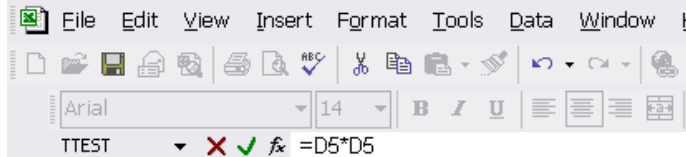
Microsoft Excel - studytime

Using Regression to explain and predict exam score from study time							
Simple Linear Regression							
	Student i	Study Time x_i	Exam Score y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	1	60	85	-48.85	0.83	-40.30	
2	2	45	80	-63.85	-4.18	266.57	
3	3	50	83	-58.85	-1.18	69.15	
4	4	75	90	-33.85	5.83	-197.18	
5	5	120	78	11.15	-6.18	-68.85	
38	38	180	80	71.15	-4.18	-297.05	
39	39	200	84	91.15	-0.17	-15.95	
40	40	210	100	101.15	15.83	1600.70	
Averages:		108.85	84.18				
Totals:							
$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$							
$\beta_0 = \bar{y} - b_1\bar{x} =$							
The regression equation ($\hat{y} = \beta_0 + \beta_1 x$):							

Copy the formula in F5 down to F44 to calculate $(x_i - \bar{x})(y_i - \bar{y})$ for the remaining observations in the data set

Calculating the coefficients of the regression line

Microsoft Excel - studytime



Calculate $(x_i - \bar{x})^2$ for the first observation in the data set:

$$G5: =D5 * D5$$

Using Regression to explain and predict							
Simple Linear Regression							
	Student i	Study Time x_i	Exam Score y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
5	1	60	85	-48.85	0.83	-40.30	=D5*D5
6	2	45	80	-63.85	-4.18	266.57	
7	3	50	83	-58.85	-1.18	69.15	
8	4	75	90	-33.85	5.83	-197.18	
9	5	120	78	11.15	-6.18	-68.85	
42	38	180	80	71.15	-4.18	-297.05	
43	39	200	84	91.15	-0.17	-15.95	
44	40	210	100	101.15	15.83	1600.70	
45	Averages:	108.85	84.18		Totals:		
47	$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$						
49	$\beta_0 = \bar{y} - b_1\bar{x} =$						
51	The regression equation ($\hat{y} = \beta_0 + \beta_1 * x$):						

Calculating the coefficients of the regression line

Microsoft Excel - studytime

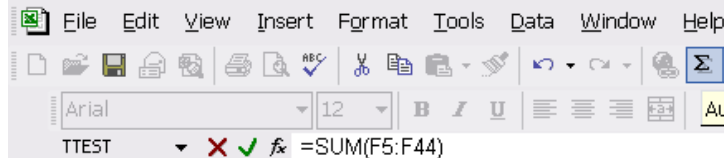
File Edit View Insert Format Tools
 Arial 14 B I
 G5 =D5*D5

Copy the formula in G5 down to G44 to calculate $(x_i - \bar{x})^2$ for the remaining observations in the data set

Using Regression to						
Simple Linear Regression						
	Study Time	Exam Score				
Student i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	60	85	-48.85	0.83	-40.30	2386.32
2	45	80	-63.85	-4.18	266.57	4076.82
3	50	83	-58.85	-1.18	69.15	3463.32
4	75	90	-33.85	5.83	-197.18	1145.82
5	120	78	11.15	-6.18	-68.85	124.323
38	180	80	71.15	-4.18	-297.05	5062.32
39	200	84	91.15	-0.17	-15.95	8308.32
40	210	100	101.15	15.83	1600.70	10231.3
Averages:	108.85	84.18		Totals:		
$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$						
$\beta_0 = \bar{y} - b_1\bar{x} =$						
The regression equation ($\hat{y} = \beta_0 + \beta_1 x$):						

Calculating the coefficients of the regression line

Microsoft Excel - studytime



Calculate $\Sigma(x_i - \bar{x})(y_i - \bar{y})$:
F45: =SUM(F5:F44)

Using Regression to explain and predict exam score from study time						
Simple Linear Regression						
	Study Time	Exam Score				
Student i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	60	85	-48.85	0.83	-40.30	2386.32
2	45	80	-63.85	-4.18	266.57	4076.82
3	50	83	-58.85	-1.18	69.15	3463.32
4	75	90	-33.85	5.83	-197.18	1145.82
5	120	78	11.15	-6.18	-68.85	124.323
38	180	80	71.15	-4.18	-297.05	5062.32
39	200	84	91.15	-0.17	-15.95	8308.32
40	210	100	101.15	15.83	1600.70	10231.3
Averages:	108.85	84.18		Totals:	=SUM(F5:F44)	
$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$						
$\beta_0 = \bar{y} - b_1\bar{x} =$						
The regression equation ($\hat{y} = \beta_0 + \beta_1 x$):						

Calculating the coefficients of the regression line

Microsoft Excel - studytime

Using Regression to explain and predict exam score from study time							
Simple Linear Regression							
	Study Time	Exam Score					
Student _i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$	
1	60	85	-48.85	0.83	-40.30	2386.32	
2	45	80	-63.85	-4.18	266.57	4076.82	
3	50	83	-58.85	-1.18	69.15	3463.32	
4	75	90	-33.85	5.83	-197.18	1145.82	
5	120	78	11.15	-6.18	-68.85	124.323	
38	180	80	71.15	-4.18	-297.05	5062.32	
39	200	84	91.15	-0.17	-15.95	8308.32	
40	210	100	101.15	15.83	1600.70	10231.3	
Averages:	108.85	84.18		Totals:	14311.05	192483.1	
$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$			$=F45/G45$				
$\beta_0 = \bar{y} - b_1\bar{x} =$							
The regression equation ($\hat{y} = \beta_0 + \beta_1^* x$):							

Calculate β_1 :

D47: =F45 / G45

Calculating the coefficients of the regression line

Microsoft Excel - studytime

File Edit View Insert Format Tools Data Window Help

Arial 14 B I U \$ % ,

TTEST X ✓ ✖ =C45-(D47*B45)

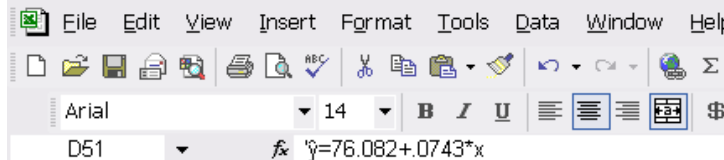
	A	B	C	D	E	F	G	H
1	Using Regression to explain and predict exam score from study time							
2	Simple Linear Regression							
3		Study Time	Exam Score					
4	Student i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$	
5	1	60	85	-48.85	0.83	-40.30	2386.32	
6	2	45	80	-63.85	-4.18	266.57	4076.82	
7	3	50	83	-58.85	-1.18	69.15	3463.32	
8	4	75	90	-33.85	5.83	-197.18	1145.82	
9	5	120	78	11.15	-6.18	-68.85	124.323	
42	38	180	80	71.15	-4.18	-297.05	5062.32	
43	39	200	84	91.15	-0.17	-15.95	8308.32	
44	40	210	100	101.15	15.83	1600.70	10231.3	
45	Averages:	108.85	84.18		Totals:	14311.05	192483.1	
47	$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$			0.0743				
49	$\beta_0 = \bar{y} - b_1\bar{x} =$			=C45-(D47*B45)				
51	The regression equation ($\hat{y} = \beta_0 + \beta_1 x$):							

Calculate β_0 :

D49: =C45 – (D47 * B45)

Calculating the coefficients of the regression line

Microsoft Excel - studytime



Write the regression equation:

$$D51: \hat{y} = 76.08 + .0743x$$

Using Regression to explain and predict exam score from study time							
Simple Linear Regression							
	Study Time	Exam Score					
Student i	x_i	y_i	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$	
1	60	85	-48.85	0.83	-40.30	2386.32	
2	45	80	-63.85	-4.18	266.57	4076.82	
3	50	83	-58.85	-1.18	69.15	3463.32	
4	75	90	-33.85	5.83	-197.18	1145.82	
5	120	78	11.15	-6.18	-68.85	124.323	
38	180	80	71.15	-4.18	-297.05	5062.32	
39	200	84	91.15	-0.17	-15.95	8308.32	
40	210	100	101.15	15.83	1600.70	10231.3	
Averages:	108.85	84.18			Totals:	14311.05	192483.1
$\beta_1 = \Sigma(x_i - \bar{x})(y_i - \bar{y}) / \Sigma(x_i - \bar{x})^2 =$			0.0743				
$\beta_0 = \bar{y} - b_1\bar{x} =$			76.082				
The regression equation ($\hat{y} = \beta_0 + \beta_1 * x$)			$\hat{y} = 76.082 + .0743 * x$				

Measuring the Fit of the Regression Model

- Regression models can be developed for any variables X and Y
- How do we know that the model is actually helpful in predicting Y based on X ?
- Three measures of variability are
 - SST – Total variability about the **mean**
 - SSE – Variability about the regression line
 - SSR – **Total variability that is explained by the model**

Fit of the Regression Model (con't)

- Sum of the squares total (actual minus average, \bar{Y} is overall average)

$$SST = \sum (Y - \bar{Y})^2$$

- Sum of the squared error (actual minus model)

$$SSE = \sum e^2 = \sum (Y - \hat{Y})^2$$

- Sum of squares due to regression (model minus overall average)

$$SSR = \sum (\hat{Y} - \bar{Y})^2$$

- An important relationship

$$SST = SSR + SSE$$

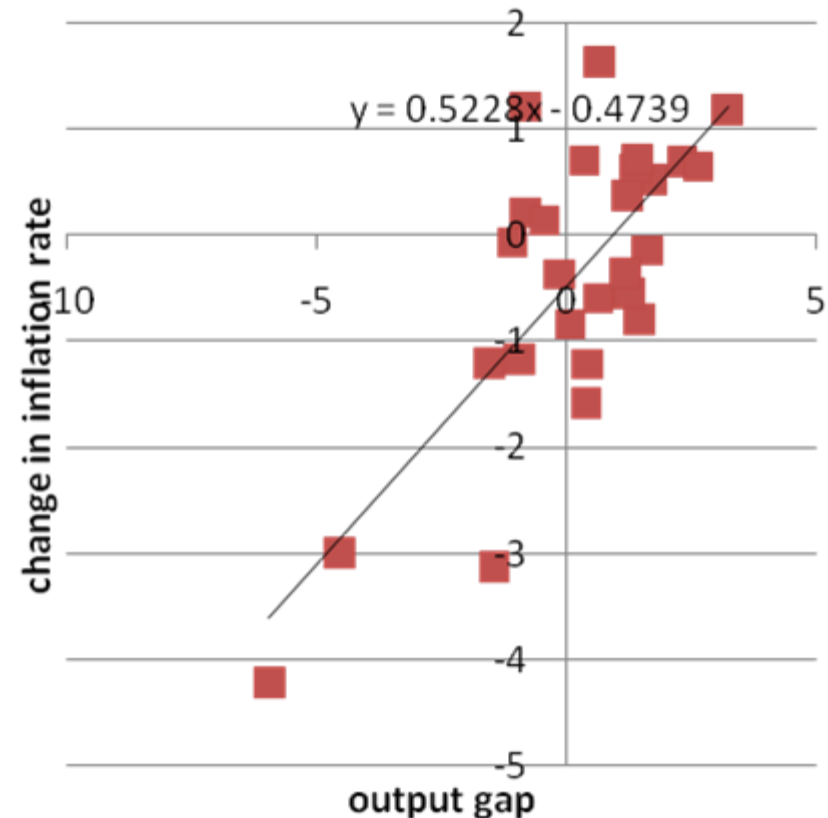
Coefficient of Determination

- The proportion of the variability in Y explained by regression equation is called the *coefficient of determination*
- The coefficient of determination is r^2 , which is between zero and 1
- A smaller SSE/SST ratio yields a larger r^2

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

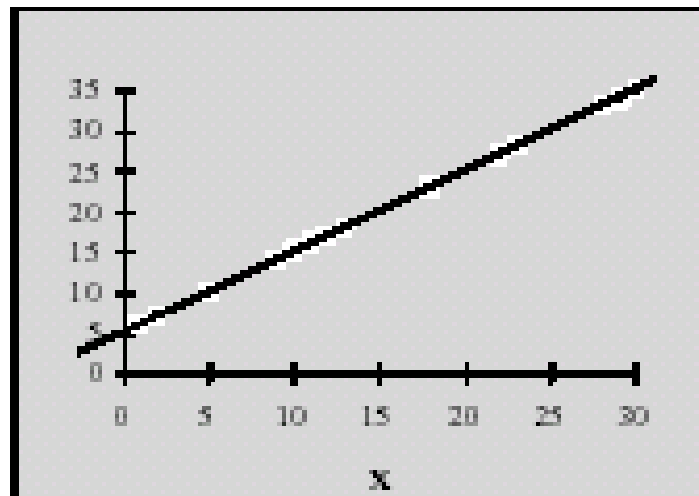
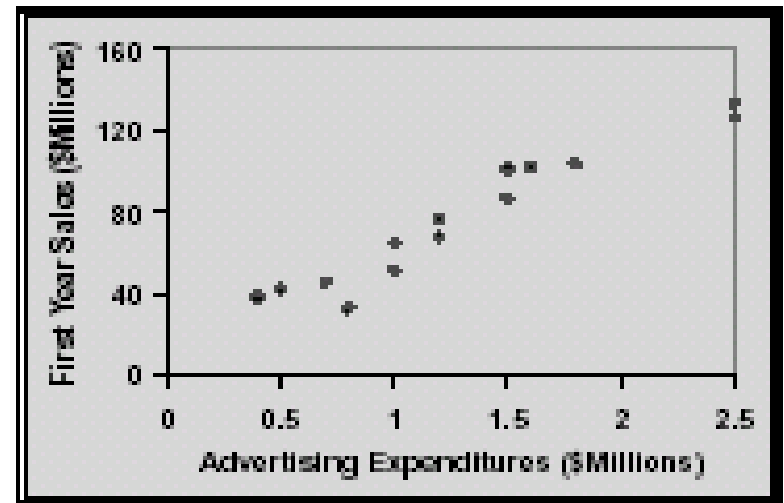
Coefficient of Determination (con't)

- It is a measure of the overall quality of the regression
- r^2 (or R^2) takes on values between 0 and 1; it is a ratio of that portion of the dependent variable's variation that is explained by the fitted model
- R^2 is pearson's r squared

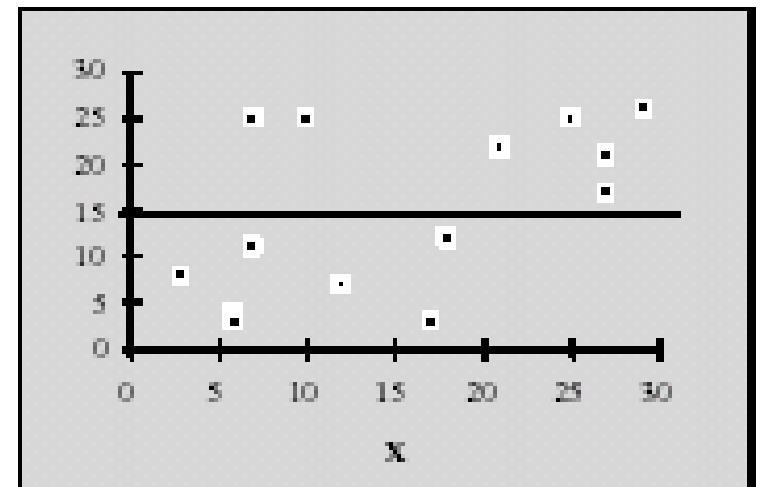


R^2 takes values between 0 and 1 (it is a percentage).

$R^2 = 0.833$ in our Appleglo Example



$R^2 = 1$; x values account for all variation in the Y values



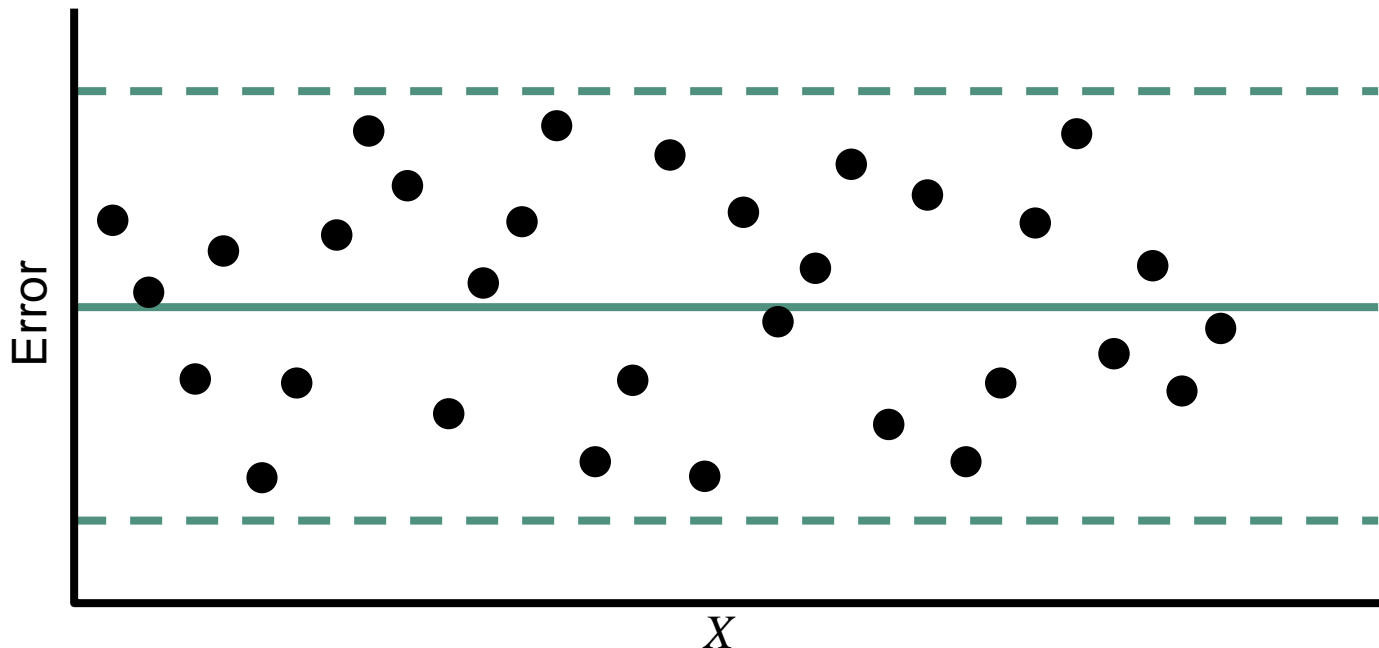
$R^2 = 0$; x values account for none variation in the Y values

Assumptions of the Regression Model

- If we make certain assumptions about the errors in a regression model, we can perform statistical tests to determine if the model is useful
 1. Errors are independent
 2. Errors are normally distributed
 3. Errors have a mean of zero
 4. Errors have a constant variance
- A plot of the residuals (errors) will often highlight any glaring violations of the assumptions

Residual Plots

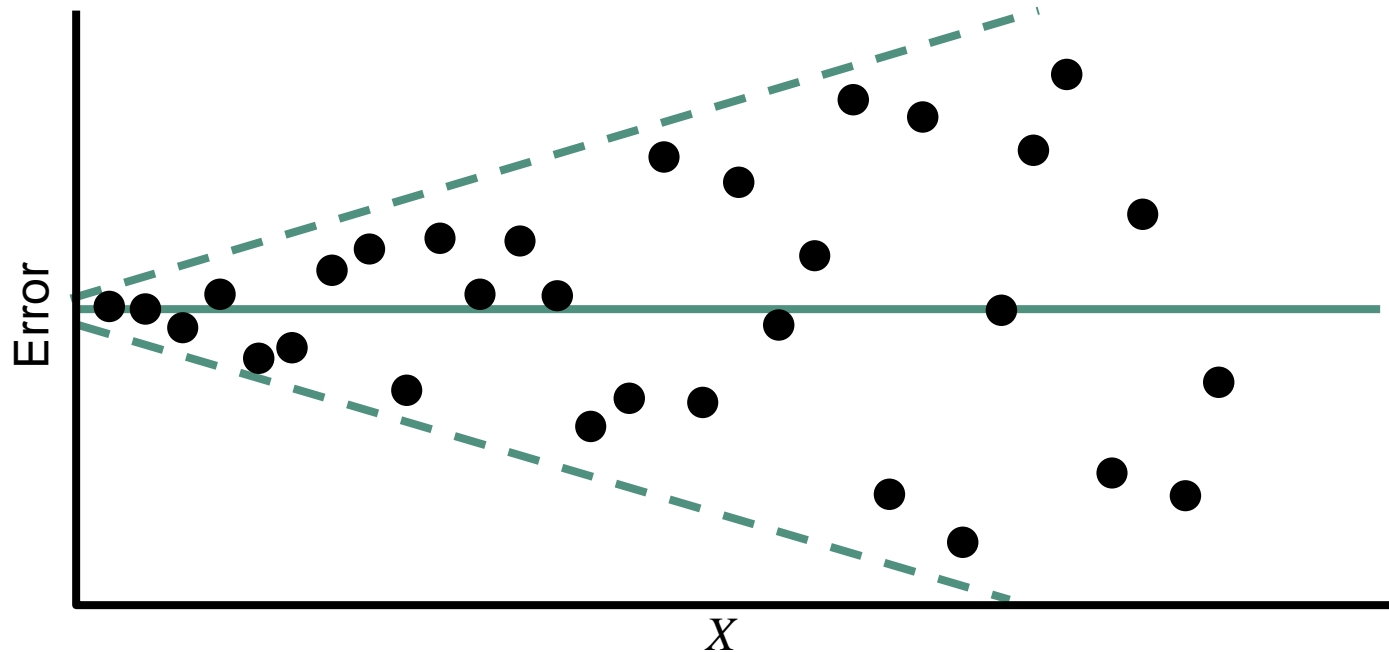
- A random plot of residuals



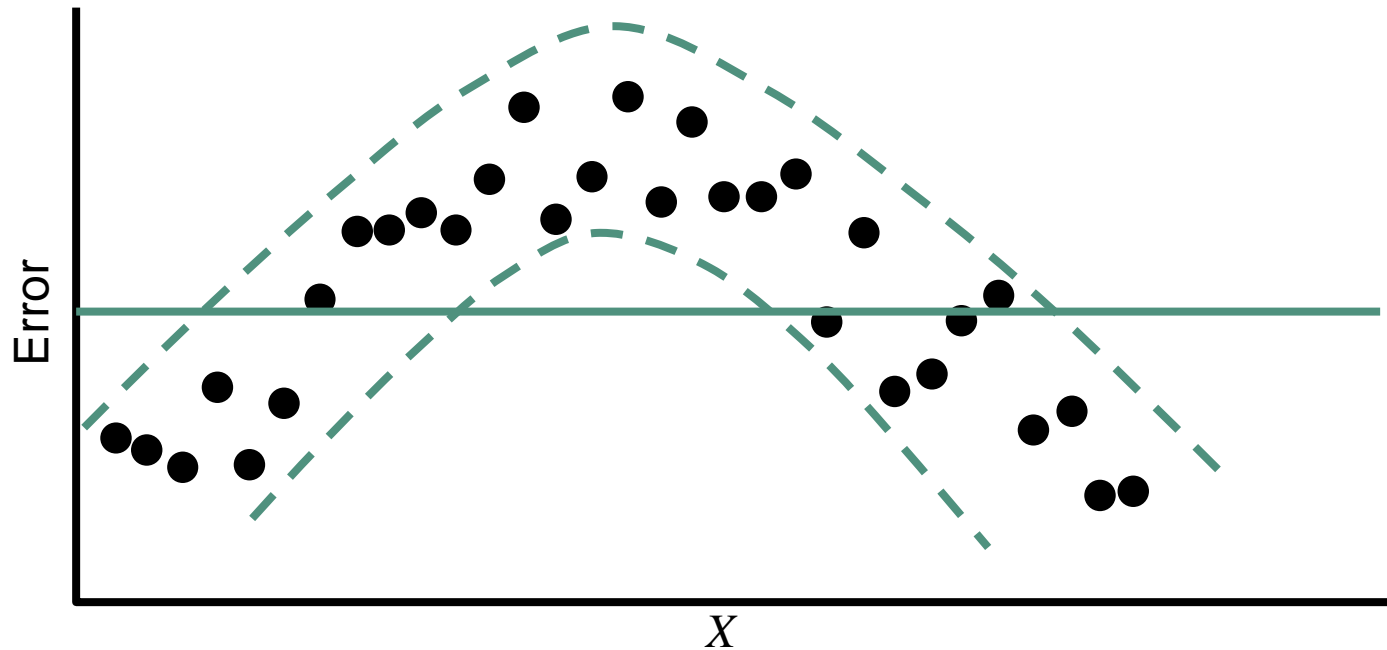
This is the type of random pattern we want to see for a good linear model.

Residual Plots (con't)

■ Nonconstant error variance



Non Linear Residual Plot



Does the western conference NBA team with the largest salary cap win more games?

Team	Win %	03-04 Salary Cap
Timberwolves	70.7	72
Lakers	68.3	66
Spurs	69.5	47
Kings	67.1	70
Mavericks	63.4	79
Grizzlies	61.0	58
Rockets	54.9	52
Nuggets	52.4	36
Jazz	51.2	28
Trailblazers	50.0	84
Warriors	45.1	52
Supersonics	45.1	51
Suns	35.4	65
Clippers	34.1	38

How is the NBA salary cap determined ?



NBA Salary Cap

- This is the limit to the total amount of money that [National Basketball Association](#) teams are allowed to pay their players
- It is defined by the [league's collective bargaining agreement](#) (CBA), and is subject to a complex system of rules and exceptions and as such is considered a "soft" cap
- The actual amount of the [salary cap](#) varies on a year-to-year basis, and is **calculated as a percentage of the league's revenue from the previous season**
- Like many professional sports leagues, the NBA has a salary cap to control cost and foster team equality
- The maximum amount of money a player can sign for is based on the number of years that player has played and the total of the salary cap; for example, the maximum salary of a player with 6 or fewer years of experience is either \$9,000,000 or 25% of the total salary cap (2013–14: \$14,670,000), whichever is greater; however there are many types of “exceptions”

VIEWING:

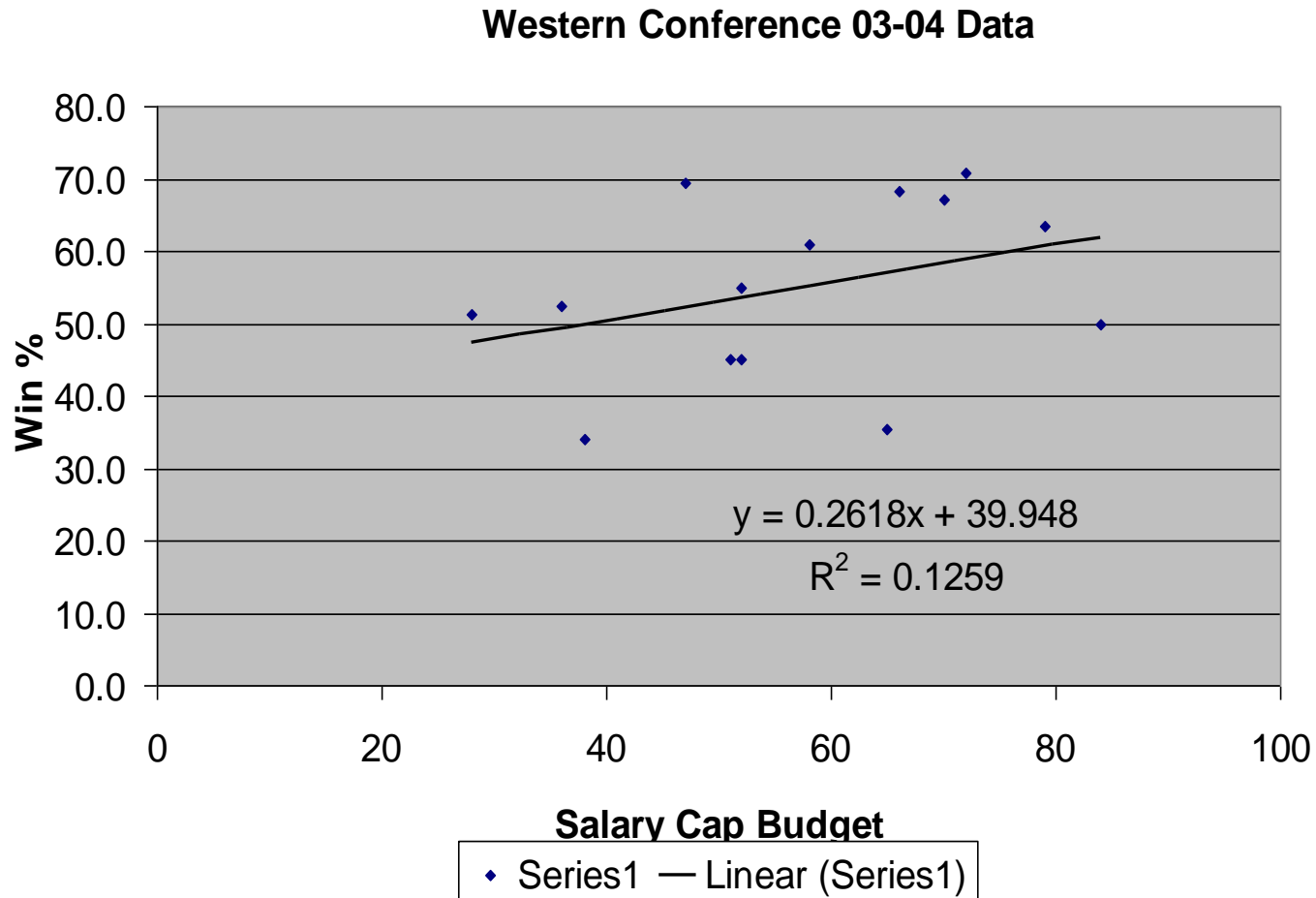
2014

Team Salary Caps

GO

RANK	TEAM	ACTIVE PLAYERS	GUARD CAP	FORWARD CAP	CENTER CAP	TOTAL CAP	MAX TYPE	SPACE
1	Atlanta Hawks	15	\$19,555,933	\$27,451,111	\$12,000,000	60,035,044	Hard	\$20,790,956
2	Philadelphia 76ers	15	\$10,141,055	\$15,023,903	\$5,730,363	42,295,914	Max	\$20,765,086
3	Phoenix Suns	15	\$35,542,326	\$12,555,460	\$5,549,920	61,541,986	Hard	\$19,157,034
4	Minnesota Timberwolves	14	\$24,093,277	\$21,946,749	\$13,561,505	64,389,797	Hard	\$16,439,203
5	Los Angeles Lakers	14	\$25,407,555	\$10,302,023	\$9,915,243	46,707,154	Max	\$16,357,546
6	Portland Trail Blazers	15	\$16,561,503	\$39,215,372	\$13,242,649	69,352,824	Hard	\$11,476,176
7	Denver Nuggets	14	\$26,554,042	\$24,234,923	\$15,394,430	69,689,395	Hard	\$11,160,605
8	Chicago Bulls	14	\$23,455,302	\$29,515,343	\$14,122,579	67,432,557	LT	\$9,395,443
9	Sacramento Kings	15	\$13,463,266	\$41,421,300	\$15,661,243	71,514,569	Hard	\$9,314,411
10	New Orleans Pelicans	14	\$37,032,120	\$20,941,925	\$10,172,212	68,780,395	LT	\$5,045,605
11	Golden State Warriors	15	\$33,260,565	\$23,450,955	\$14,902,335	72,965,135	Hard	\$7,540,572
12	Milwaukee Bucks	16	\$19,767,155	\$19,060,161	\$16,200,000	55,997,659	Max	\$7,167,331
13	San Antonio Spurs	16	\$32,264,116	\$19,253,559	\$17,250,000	69,961,359	LT	\$6,547,541
14	Indiana Pacers	15	\$16,145,406	\$39,247,262	\$19,406,374	74,795,942	Hard	\$6,030,056
15	Memphis Grizzlies	15	\$27,632,596	\$27,919,671	\$15,529,655	75,019,956	Hard	\$5,509,042
16	Houston Rockets	15	\$26,100,274	\$25,159,615	\$21,436,271	75,771,892	Hard	\$5,057,105
17	Washington Wizards	14	\$24,557,364	\$24,221,164	\$26,350,025	75,529,553	Hard	\$5,000,447
18	Miami Heat	14	\$22,276,546	\$26,935,752	\$22,230,763	71,961,393	LT	\$4,561,607
19	Toronto Raptors	15	\$36,743,960	\$32,324,671	\$6,356,253	76,096,323	Hard	\$4,732,477
20	Orlando Magic	15	\$17,555,450	\$21,002,006	\$2,751,260	58,530,516	Max	\$4,234,452
21	Utah Jazz	17	\$26,350,200	\$24,751,992	\$6,522,074	59,582,424	Max	\$3,452,576
22	Dallas Mavericks	14	\$30,673,404	\$25,993,355	\$15,725,050	74,385,967	LT	\$2,480,433
23	Charlotte Hornets	14	\$23,962,206	\$15,456,403	\$21,404,455	60,523,067	Max	\$2,241,933
24	Los Angeles Clippers	13	\$35,560,663	\$23,570,430	\$16,745,123	79,344,759	Hard	\$1,554,245
25	Boston Celtics	14	\$21,570,136	\$31,559,465	\$1,703,760	62,271,547	Max	\$793,153
26	Detroit Pistons	14	\$21,562,460	\$14,595,656	\$11,545,293	63,034,444	Max	\$30,556
27	Oklahoma City Thunder	15	\$25,245,837	\$35,094,229	\$11,539,302	76,422,492	LT	\$-1,593,492
28	New York Knicks	15	\$11,552,634	\$62,654,950	\$915,243	80,357,678	LT	\$-3,525,675
29	Cleveland Cavaliers	15	\$15,555,390	\$45,949,354	\$16,565,233	81,593,356	LT	\$-5,064,356
30	Brooklyn Nets	14	\$52,157,649	\$7,153,436	\$29,592,625	90,549,514	LT	\$-13,720,514

Salary Cap vs Win %



Analysis of Trendline

- The prediction or “trendline” formula is:
 - $\text{Wins} = 0.2618 * \text{Salary} + 39.948$
 - So for about every 10 million of salary cap, the team will win 2.6% more games

The teams would win about 40% of their games if the players played for free ???

Determination Coefficient Analysis

■ $R^2 = 0.1259$

- This is the measure of the overall quality of the linear regression
- This value is between 0 and 1
- Here it is closer to zero, so as to suggest that there is not a very good correlation between salary and win%

Multiple Linear Regression

- SLR only involves one independent variable
- When there is more than one independent variable, then the analysis becomes more complicated
- The concepts are the same, except we are working in multi-dimensional space instead of just two dimensions
 - $Y = b_0 + b_1 * X_1 + b_2 * X_2 + \dots + b_n * X_n$
- For our advertising example here, there may be other factors influencing sales values in addition to advertising expenditures

Multiple Independent Variables

region	sales	advertising	promotions	competitor's sales
Selkirk	101.8	1.3	0.2	20.40
Susquehanna	44.4	0.7	0.2	30.50
Kittery	108.3	1.4	0.3	24.60
Acton	85.1	0.5	0.4	19.60
Finger Lakes	77.1	0.5	0.6	25.50
Berkshire	158.7	1.9	0.4	21.70
Central	180.4	1.2	1.0	6.80
Providence	64.2	0.4	0.4	12.60
Nashua	74.6	0.6	0.5	31.30
Dunster	143.4	1.3	0.6	18.60
Endicott	120.6	1.6	0.8	19.90
Five-Towns	69.7	1.0	0.3	25.60
Waldeboro	67.8	0.8	0.2	27.40
Jackson	106.7	0.6	0.5	24.30
Stowe	119.6	1.1	0.3	13.70

What are the independent and dependent variables ?

Minimize the Sum of the Residual Squares

- Choosing the coefficients ($b_0 \dots b_n$) to minimize the sum of the residual squares gives us the prediction formula
- The prediction formula here is:
 - $\text{Sales} = 65.705 + (48.979 * \text{Advertising}) + (59.654 * \text{Promotions}) - (1.838 * \text{CompetitorSales})$





Typical MLR Tool Output

Regression Statistics

Multiple R	0.913
R Square	0.833
Adjusted R Square	0.787
Standard Error	17.600
Observations	15

Analysis of Variance

	df	Sum of Squares	Mean Square	F	Significance F
Regression	3	16997.537	5665.85	18.290	0.000
Residual	11	3407.473	309.77		
Total	14	20405.009			

	Coefficients	Standard Error	t Statistic	P-value	Lower 95%	Upper 95%
Intercept	65.71	27.73	2.37	0.033	4.67	126.74
Advertising	48.98	10.66	4.60	0.000	25.52	72.44
Promotions	59.65	23.63	2.53	0.024	7.66	111.65
Competitor's Sales	-1.84	0.81	-2.26	0.040	-3.63	-0.047

Typical MLR Tool Output (con't)

$b_0 = 65.705$ (its interpretation is context dependent .

$b_1 = 48.979$ (an additional \$1 million in advertising is expected to result in an additional \$49 million in sales)

$b_2 = 59.654$ (an additional \$1 million in promotions is expected to result in an additional \$60 million in sales)

$b_3 = -1.838$ (an increase of \$1 million in competitor sales is expected to decrease sales by \$1.8 million)



Adjusted R Squared

- The formula for R squared is:
 - $R^2 = 1 - (\text{variation not accounted for} / \text{total variation})$
 - $R^2 = 1 - (SS_{\text{error}} / SS_{\text{total}})$
- The Adjusted R Squared adjusts for the degrees of freedom in the model, **and penalizes an unnecessarily complex model (too many independent variables)**
- The formula is
 - $aR^2 = 1 - (SS_{\text{error}} / df_{\text{error}}) / (SS_{\text{total}} / df_{\text{total}})$
 - Where df_{error} is the degree of freedom $[n-1-k]$
 - And df_{total} is total degrees of freedom $[n-1]$

Testing the Model for Significance

- When the sample size is too small, you can get good values for MSE and r^2 even if there is no relationship between the variables
- Testing the model for significance helps determine if the values are meaningful
- We do this by performing a statistical hypothesis test



Testing the Model for Significance (con't)

- The overall significance is checked by performing an F-test, with the null hypothesis being that all the slopes are zero
 - When F is large then the significance is small and the overall model is good
- We can also do a hypothesis test for each slope (variable) using a t-test
 - When the p-value is small, then there is a strong relationship for that variable
- For simple linear regression, the F test and t-test give the same results

Testing the Model for Significance (con't)

- The F statistic is based on the MSE and MSR

$$MSR = \frac{SSR}{k}$$

where

k = number of independent variables in the model

- The F statistic is

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

- This describes an F distribution with
degrees of freedom for the numerator = $df_1 = k$
degrees of freedom for the denominator = $df_2 = n - k - 1$

Testing the Model for Significance (con't)

Regression Statistics

Multiple R	0.913
R Square	0.833
Adjusted R Square	0.787
Standard Error	17.600
Observations	15

Analysis of Variance

	df	Sum of Squares	Mean Square	F	Significance F
Regression	3	16997.537	5665.85	18.290	0.000
Residual	11	3407.473	309.77		
Total	14	20405.009			

	Coefficients	Standard Error	t Statistic	P-value	Lower 95%	Upper 95%
Intercept	65.71	27.73	2.37	0.033	4.67	126.74
Advertising	48.98	10.66	4.60	0.000	25.52	72.44
Promotions	59.65	23.63	2.53	0.024	7.66	111.65
Competitor's Sales	-1.84	0.81	-2.26	0.040	-3.63	-0.047

Testing the Model for Significance (con't)

- Looking at the coefficients (slopes), can we tell which independent variables are the most influential ?

Regression Statistics							
Multiple R				0.913			
R Square				0.833			
Adjusted R Square				0.787			
Standard Error				17.600			
Observations				15			
Analysis of Variance							
	df	Sum of Squares	Mean Square	F	Significance F		
Regression	3	16997.537	5665.85	18.290	0.000		
Residual	11	3407.473	309.77				
Total	14	20405.009					
	Coefficients	Standard Error	t Statistic	P-value	Lower 95%	Upper 95%	
Intercept	65.71	27.73	2.37	0.033	4.67	126.74	
Advertising	48.98	10.66	4.60	0.000	25.52	72.44	
Promotions	59.65	23.63	2.53	0.024	7.66	111.65	
Competitor's Sales	-1.84	0.81	-2.26	0.040	-3.63	-0.047	

Multicollinearity

- **Multicollinearity** is a statistical phenomenon in which two or more predictor variables in a multiple regression model are highly correlated
- In this situation the coefficient estimates may change erratically in response to small changes in the model or the data
- Multicollinearity does not reduce the predictive power or reliability of the model as a whole – the overall F test is still valid
- It only affects calculations regarding individual predictors (independent variables)

Singularity

- Singularity is the extreme form of [multicollinearity](#) - when a near perfect linear relationship exists between variables
- Such absolute multicollinearity could arise when independent variables are linearly related in their definition
- A simple example: two variables "height in centimeters" and "height in inches" are included in the regression model

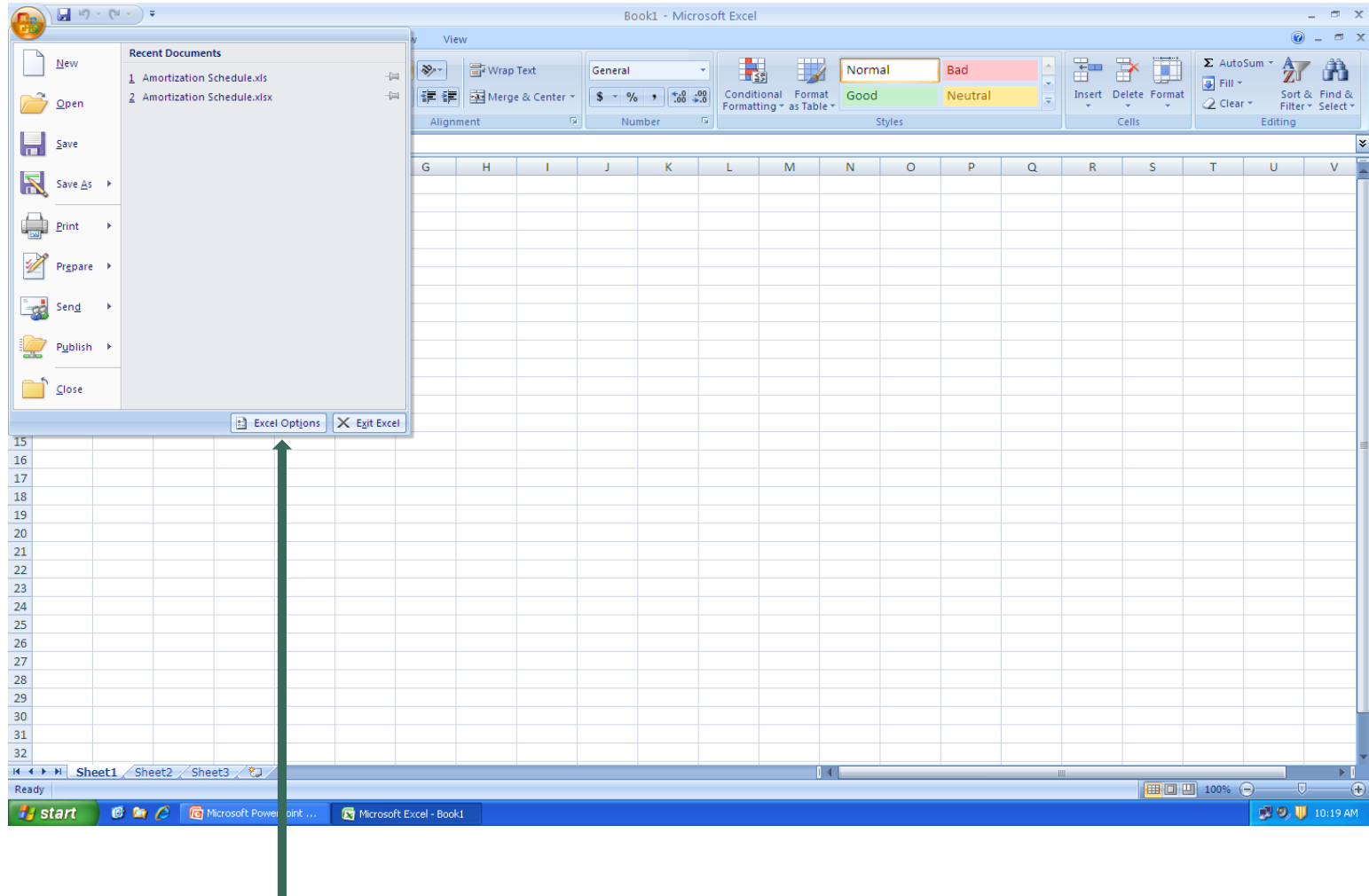
Matrix Solution of MLR

- To perform MLR, a matrix formulation and solution is normally used
- $b = (X'X)^{-1} X'y$
 - where X is a matrix of the data values for the independent variables, y is a vector of the values of the dependent variable, b is the solution vector (weighting parameters), $^{-1}$ denotes the matrix inverse, and $'$ denotes the transpose of the matrix.
- Numerically it is not efficient for large models to directly calculate matrix inverses, so other numerical techniques are used such as Gauss, Householder, etc.
- “Sparse Matrix” techniques are also used for large regression problems

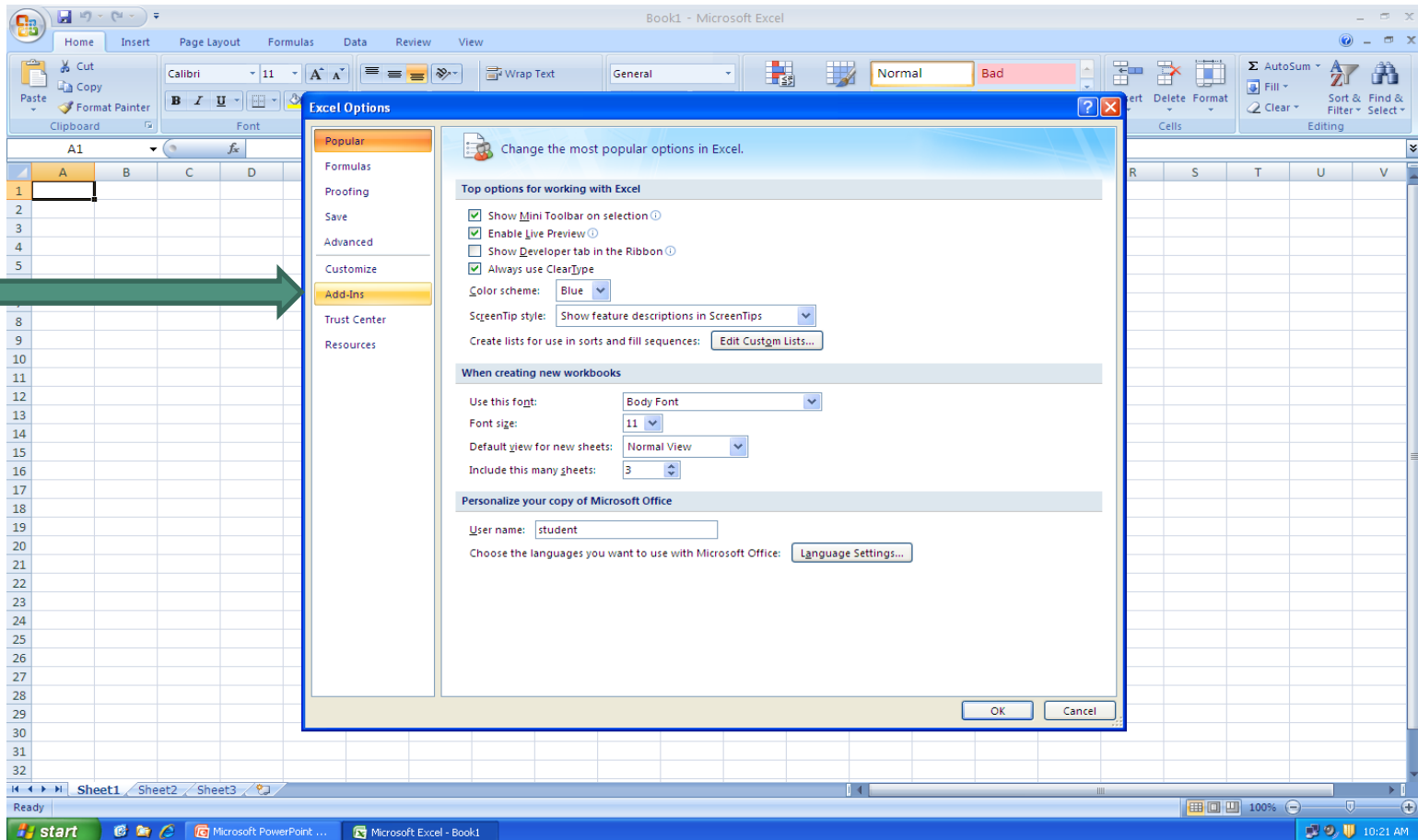
MLR in Excel

- Install the “Analysis ToolPak” if not already installed
 - Excel 2003: Tools...Add-ins...Analysis ToolPak
 - Excel 2007/2010+: Office Button (File in 2010)...Excel Options...Add-ins...Analysis ToolPak
- Enter dependent and independent variable data (see next slide)
- Tools...Data Analysis...Regression

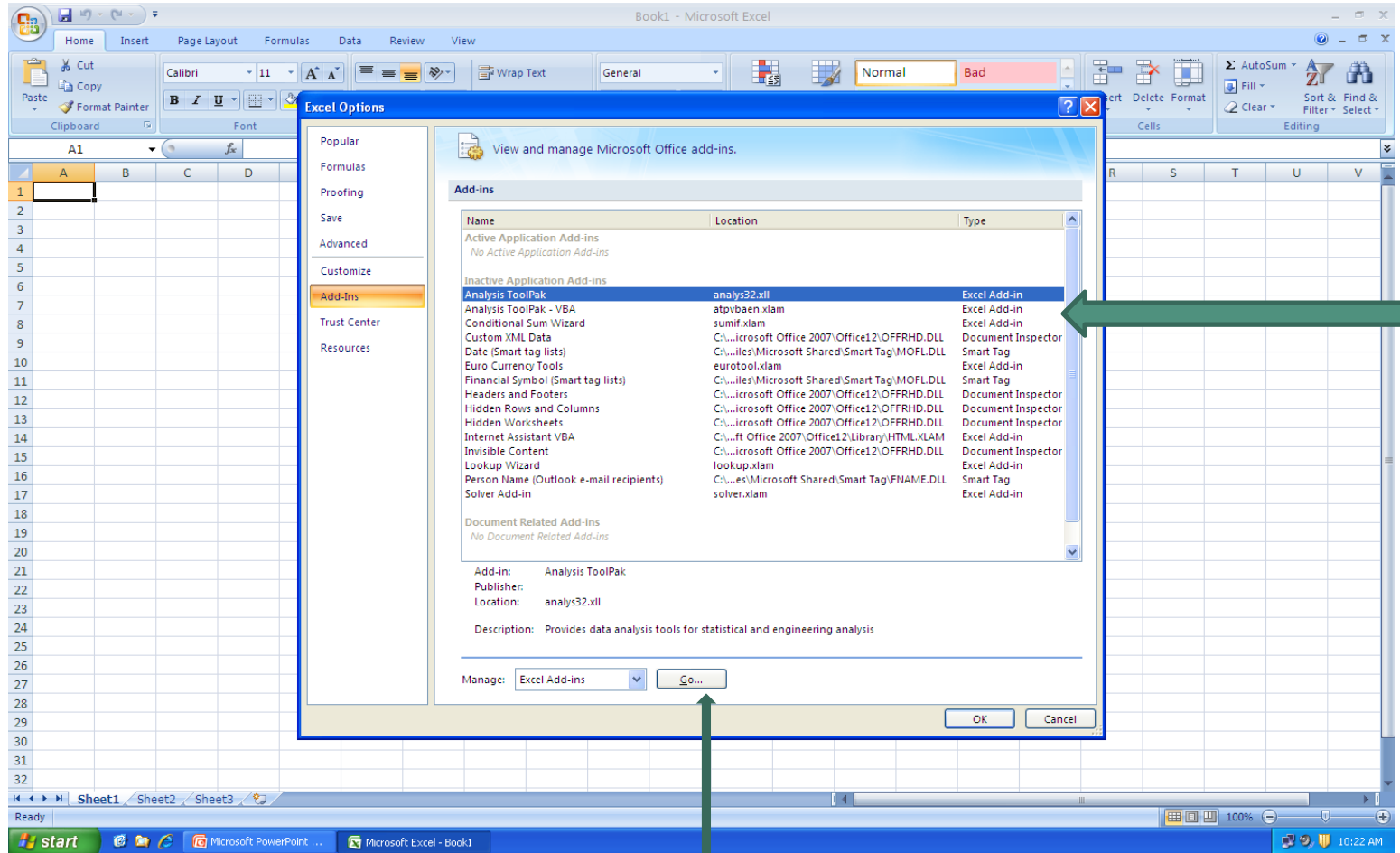
Office Button → Excel Options



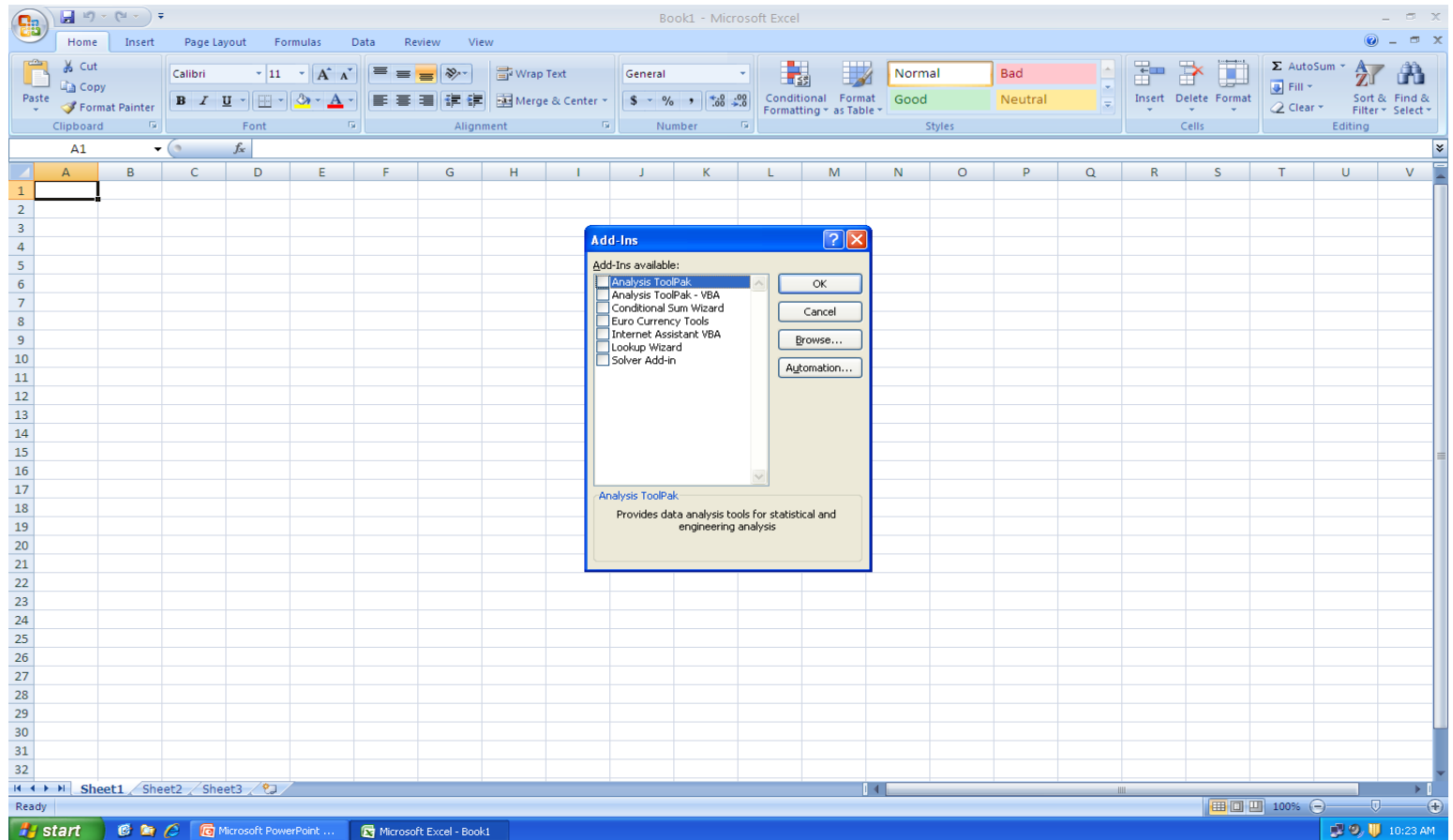
Add-Ins



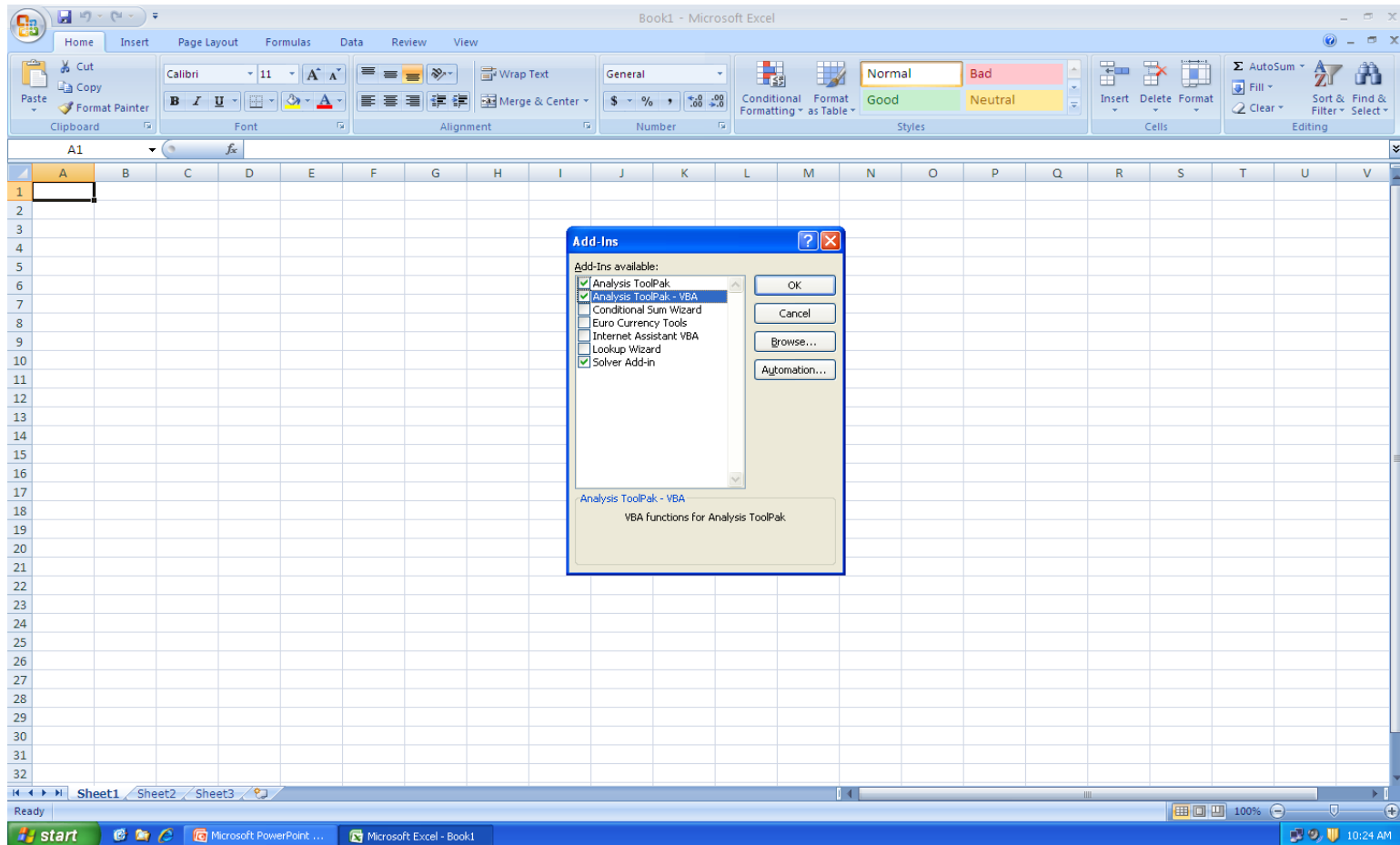
Analysis ToolPak → Go



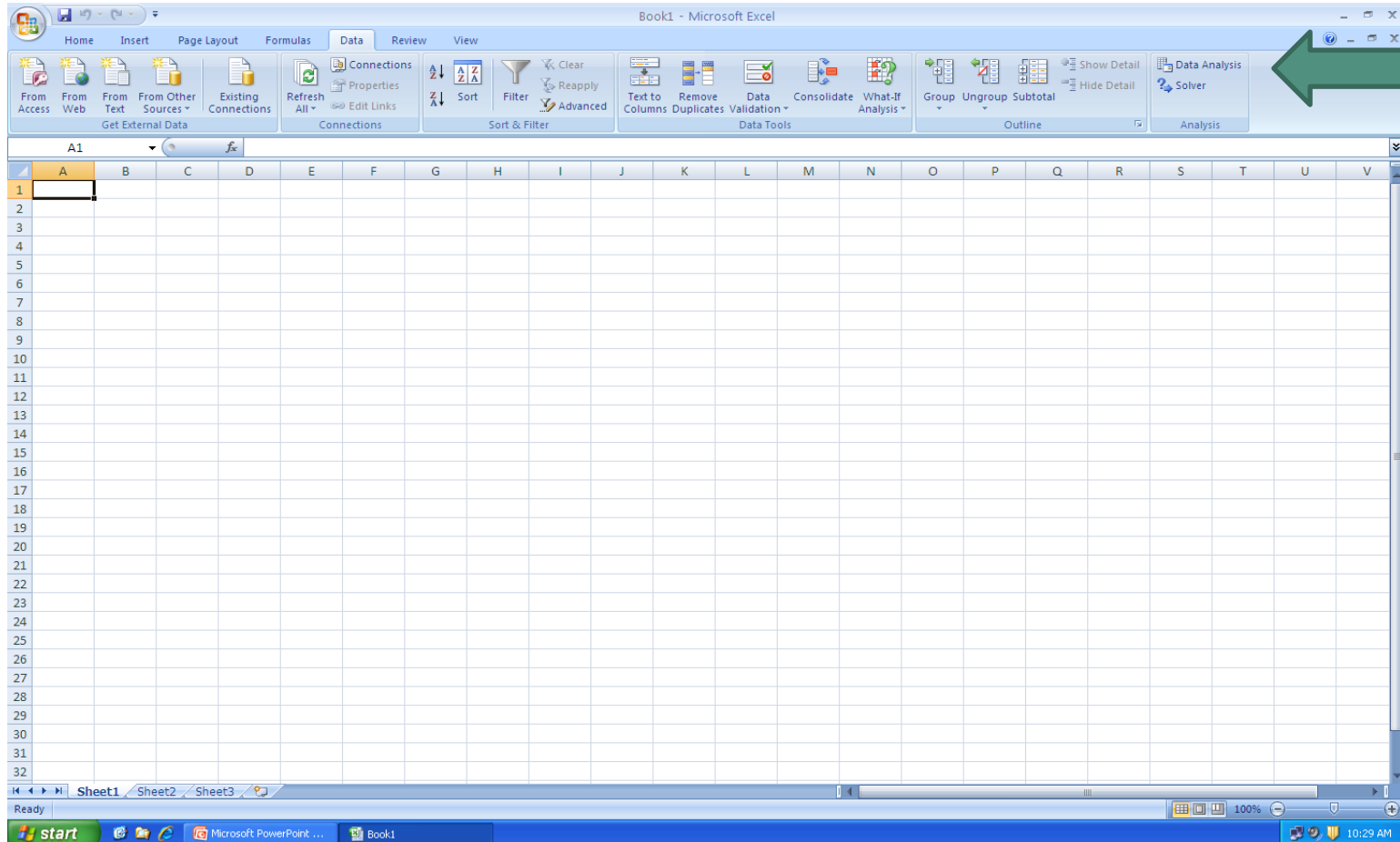
Add In Selections



Check-off Options → Ok



Data Tab → Analysis



1



Regression in Excel (con't)



StudyTime.xls - Microsoft Excel

Home Insert Page Layout Formulas Data Review View Developer

Get External Data Refresh All Properties Edit Links Connections Sort & Filter Sort Filter Clear Reapply Advanced Text to Columns Remove Duplicates Data Tools Data Validation Consolidate What-If Analysis Group Ungroup Subtotal Outline Analysis Solver

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	Student	StudyTimeMin	ExamScore											
2	1	60	85											
3	2	45	80											
4	3	50	83											
5	4	75	90											
6	5	120	78											
7	6	180	94											
8	7	240	90											
9	8	45	84											
10	9	30	70											
11	10	90	93											
12	11	120	90											
13	12	180	90											
14	13	210	90											
15	14	240	92											
16	15	90	89											
17	16	75	77											
18	17	90	80											
19	18	120	85											
20	19	240	95											
21	20	270	94											
22	21	15	68											
23	22	30	75											
24	23	45	85											

Regression

Input

Input Y Range: \$C\$1:\$C\$41

Input X Range: \$B\$1:\$B\$41

☒ Labels ☐ Constant is Zero

☒ Confidence Level: 95 %

Output options

☐ Output Range:

☒ New Worksheet Ply:

☐ New Workbook

Residuals

☐ Residuals ☐ Residual Plots

☐ Standardized Residuals ☐ Line Fit Plots

Normal Probability

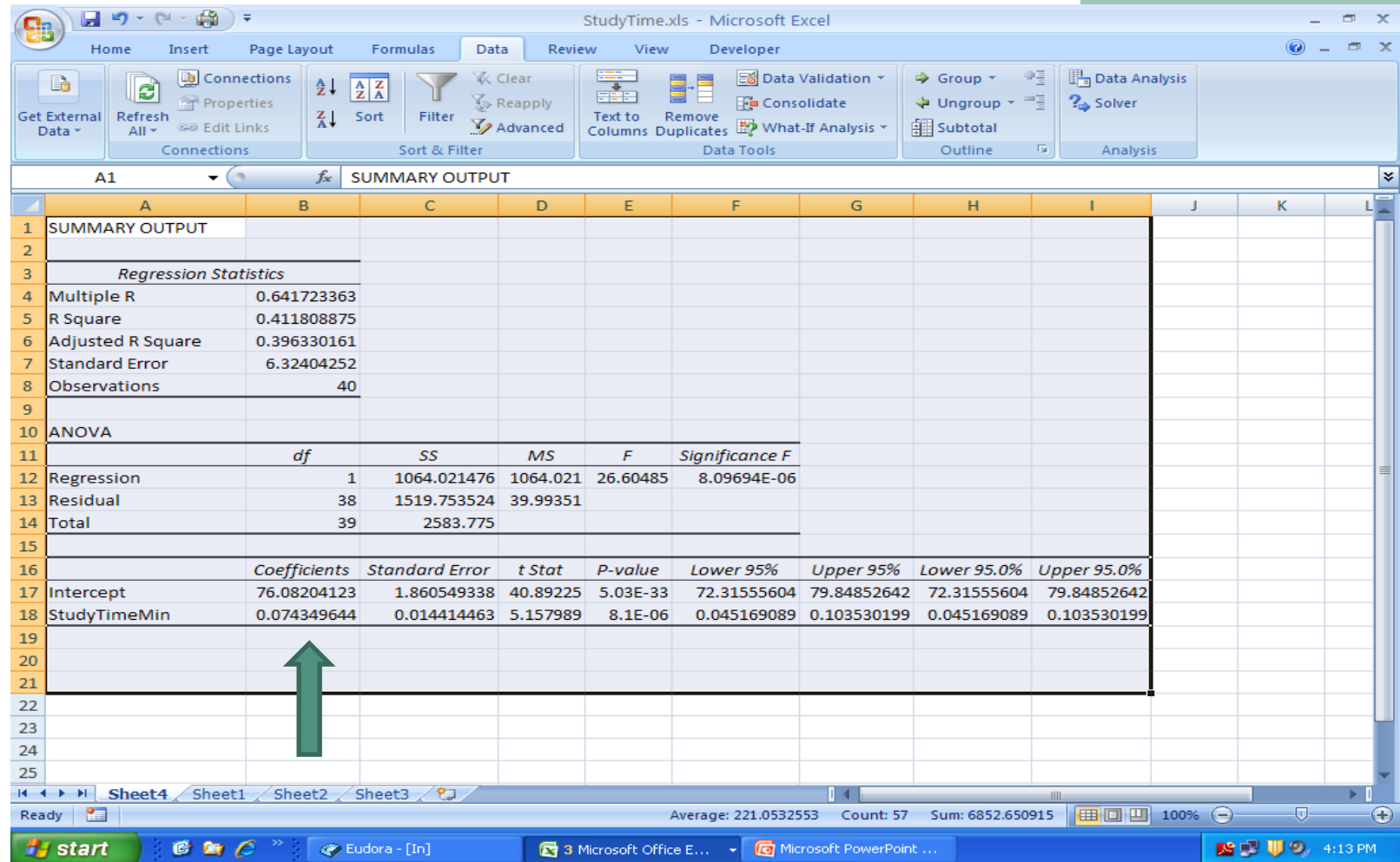
☐ Normal Probability Plots

Sheet1 Sheet2 Sheet3

Ready 100%

start Eudora - [In] 3 Microsoft Office E... Microsoft PowerPoint ... 4:12 PM

Coefficients: Intercept & Slope



StudyTime.xls - Microsoft Excel

Home Insert Page Layout Formulas Data Review View Developer

Get External Data Refresh All Connections Sort Filter Clear Reapply Advanced Text to Columns Remove Duplicates Data Validation Consolidate What-If Analysis Group Ungroup Subtotal Outline Data Analysis Solver

A1 SUMMARY OUTPUT

	A	B	C	D	E	F	G	H	I	J	K	L
1	SUMMARY OUTPUT											
2												
3	Regression Statistics											
4	Multiple R	0.641723363										
5	R Square	0.411808875										
6	Adjusted R Square	0.396330161										
7	Standard Error	6.32404252										
8	Observations	40										
9												
10	ANOVA											
11		df	SS	MS	F	Significance F						
12	Regression	1	1064.021476	1064.021	26.60485	8.09694E-06						
13	Residual	38	1519.753524	39.99351								
14	Total	39	2583.775									
15												
16		Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%			
17	Intercept	76.08204123	1.860549338	40.89225	5.03E-33	72.31555604	79.84852642	72.31555604	79.84852642			
18	StudyTimeMin	0.074349644	0.014414463	5.157989	8.1E-06	0.045169089	0.103530199	0.045169089	0.103530199			
19												
20												
21												
22												
23												
24												
25												

Sheet4 Sheet1 Sheet2 Sheet3

Ready Average: 221.0532553 Count: 57 Sum: 6852.650915 100%

start Eudora - [In] 3 Microsoft Office E... Microsoft PowerPoint ... 4:13 PM

Regression in Excel (con't)

[F and p-value are the same for SLR]

StudyTime.xls - Microsoft Excel

Home Insert Page Layout Formulas Data Review View Developer

Get External Data Refresh All Connections Properties Edit Links Connections Sort & Filter Filter Reapply Advanced Text to Columns Remove Duplicates Data Tools Data Validation Consolidate What-If Analysis Group Ungroup Subtotal Outline Analysis Data Analysis Solver

A1 SUMMARY OUTPUT

	A	B	C	D	E	F	G	H	I	J	K	L
1	SUMMARY OUTPUT											
2												
3	Regression Statistics											
4	Multiple R	0.641723363										
5	R Square	0.411808875										
6	Adjusted R Square	0.396330161										
7	Standard Error	6.32404252										
8	Observations	40										
9												
10	ANOVA											
11		df	SS	MS	F	Significance F						
12	Regression	1	1064.021476	1064.021	26.60485	8.09694E-06						
13	Residual	38	1519.753524	39.99351								
14	Total	39	2583.775									
15												
16		Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%			
17	Intercept	76.08204123	1.860549338	40.89225	5.03E-33	72.31555604	79.84852642	72.31555604	79.84852642			
18	StudyTimeMin	0.074349644	0.014414463	5.157989	8.1E-06	0.045169089	0.103530199	0.045169089	0.103530199			
19												
20												
21												
22												
23												
24												
25												

Sheet4 Sheet1 Sheet2 Sheet3

Ready Average: 221.0532553 Count: 57 Sum: 6852.650915 100%

start Eudora - [In] 3 Microsoft Office E... Microsoft PowerPoint ... 4:13 PM

Dimensions of Regression Coefficients

- The units of the intercept are the same units as the dependent variable
- The units of the coefficients for the independent variables are the dependent variable units divided by the independent variable units
- For example in the relationship between house price in dollars (dependent variable) and house square footage (independent variable):
 - The intercept is in dollars
 - The slope is in dollars per square foot

- For our example here of study time (in minutes) versus exam score (in %):
 - What are the units of the intercept ?
 - What are the units of the slope ?



Wait....



Don't look ahead, until
you have your answer !

- For our example here of study time (in minutes) versus exam score (points):
 - Units of the intercept - points
 - Units of the slope - points/min



Using The Regression Equation For Prediction

- If you wanted to predict what your exam score would be if you studied for 90 minutes:

$$\hat{y} = 76.072 + .0743(90)$$

$$\hat{y} = 82.7$$

- You can expect to get an 82.7% on the exam if you study for 90 minutes
- The regression equation is supposed to improve our ability to predict an outcome value (exam score) because we know the value of the related variable (study time)

MLR in Excel - Sales Data Set

region	sales	advertising	promotions	competitor's sales
Selkirk	101.8	1.3	0.2	20.40
Susquehanna	44.4	0.7	0.2	30.50
Kittery	108.3	1.4	0.3	24.60
Acton	85.1	0.5	0.4	19.60
Finger Lakes	77.1	0.5	0.6	25.50
Berkshire	158.7	1.9	0.4	21.70
Central	180.4	1.2	1.0	6.80
Providence	64.2	0.4	0.4	12.60
Nashua	74.6	0.6	0.5	31.30
Dunster	143.4	1.3	0.6	18.60
Endicott	120.6	1.6	0.8	19.90
Five-Towns	69.7	1.0	0.3	25.60
Waldeboro	67.8	0.8	0.2	27.40
Jackson	106.7	0.6	0.5	24.30
Stowe	119.6	1.1	0.3	13.70



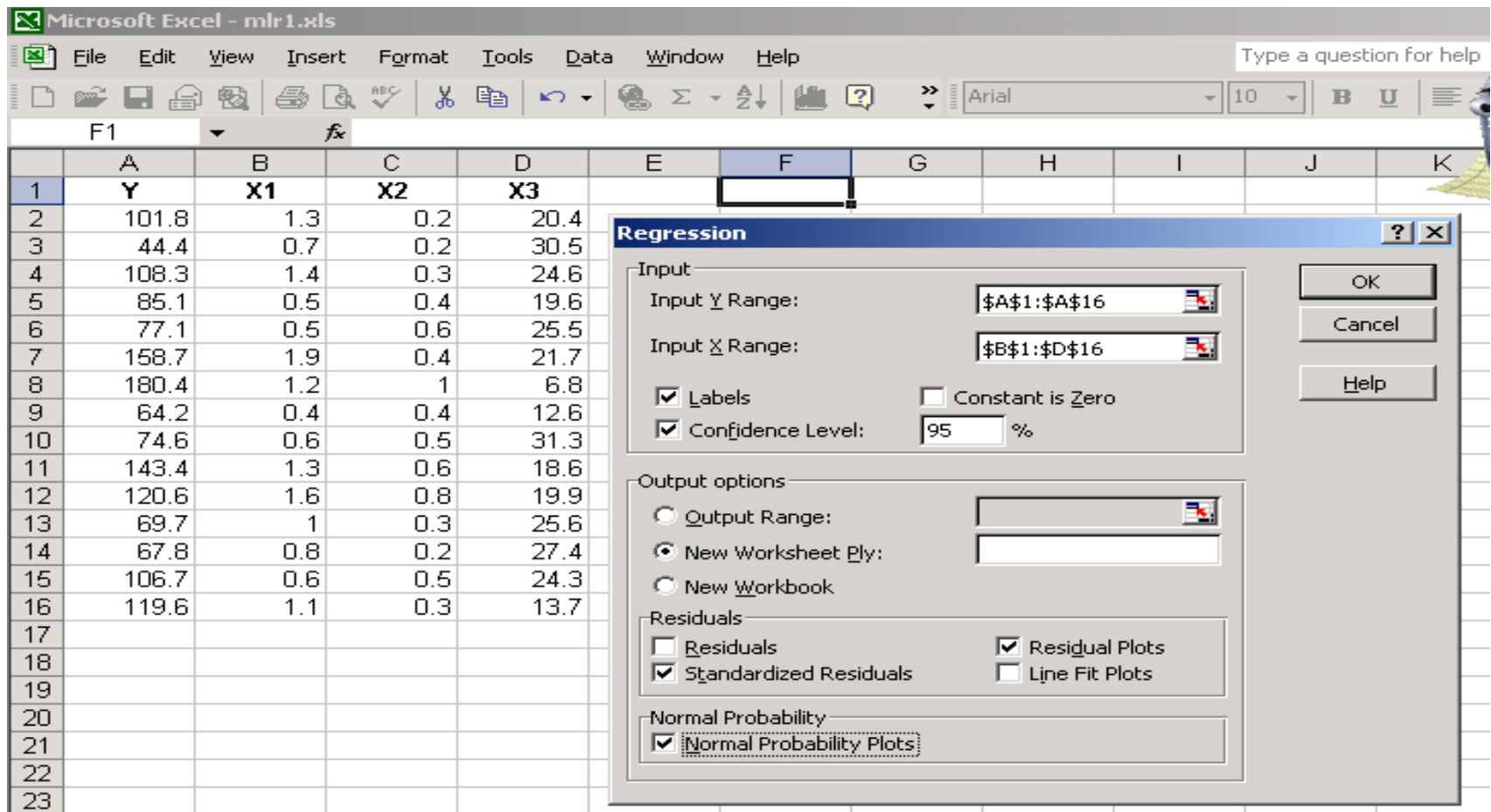
Microsoft Excel - mlr1.xls

File Edit View Insert Format Tools Data

F1 fx

	A	B	C	D
1	Y	X1	X2	X3
2	101.8	1.3	0.2	20.4
3	44.4	0.7	0.2	30.5
4	108.3	1.4	0.3	24.6
5	85.1	0.5	0.4	19.6
6	77.1	0.5	0.6	25.5
7	158.7	1.9	0.4	21.7
8	180.4	1.2	1	6.8
9	64.2	0.4	0.4	12.6
10	74.6	0.6	0.5	31.3
11	143.4	1.3	0.6	18.6
12	120.6	1.6	0.8	19.9
13	69.7	1	0.3	25.6
14	67.8	0.8	0.2	27.4
15	106.7	0.6	0.5	24.3
16	119.6	1.1	0.3	13.7
17				
18				

Regression



The screenshot shows the Microsoft Excel interface with a workbook named 'mlr1.xls'. The data is organized in columns A through D, with row 1 containing labels (Y, X1, X2, X3) and rows 2 through 16 containing numerical data. The 'Regression' dialog box is open, showing the following settings:

- Input:**
 - Input Y Range: \$A\$1:\$A\$16
 - Input X Range: \$B\$1:\$D\$16
 - ☒ Labels
 - ☐ Constant is Zero
 - ☒ Confidence Level: 95 %
- Output options:**
 - ☐ Output Range:
 - ☒ New Worksheet Ply:
 - ☐ New Workbook
- Residuals:**
 - ☐ Residuals
 - ☒ Standardized Residuals
 - ☒ Residual Plots
 - ☐ Line Fit Plots
- Normal Probability:**
 - ☒ Normal Probability Plots

	A	B	C	D
1	Y	X1	X2	X3
2	101.8	1.3	0.2	20.4
3	44.4	0.7	0.2	30.5
4	108.3	1.4	0.3	24.6
5	85.1	0.5	0.4	19.6
6	77.1	0.5	0.6	25.5
7	158.7	1.9	0.4	21.7
8	180.4	1.2	1	6.8
9	64.2	0.4	0.4	12.6
10	74.6	0.6	0.5	31.3
11	143.4	1.3	0.6	18.6
12	120.6	1.6	0.8	19.9
13	69.7	1	0.3	25.6
14	67.8	0.8	0.2	27.4
15	106.7	0.6	0.5	24.3
16	119.6	1.1	0.3	13.7

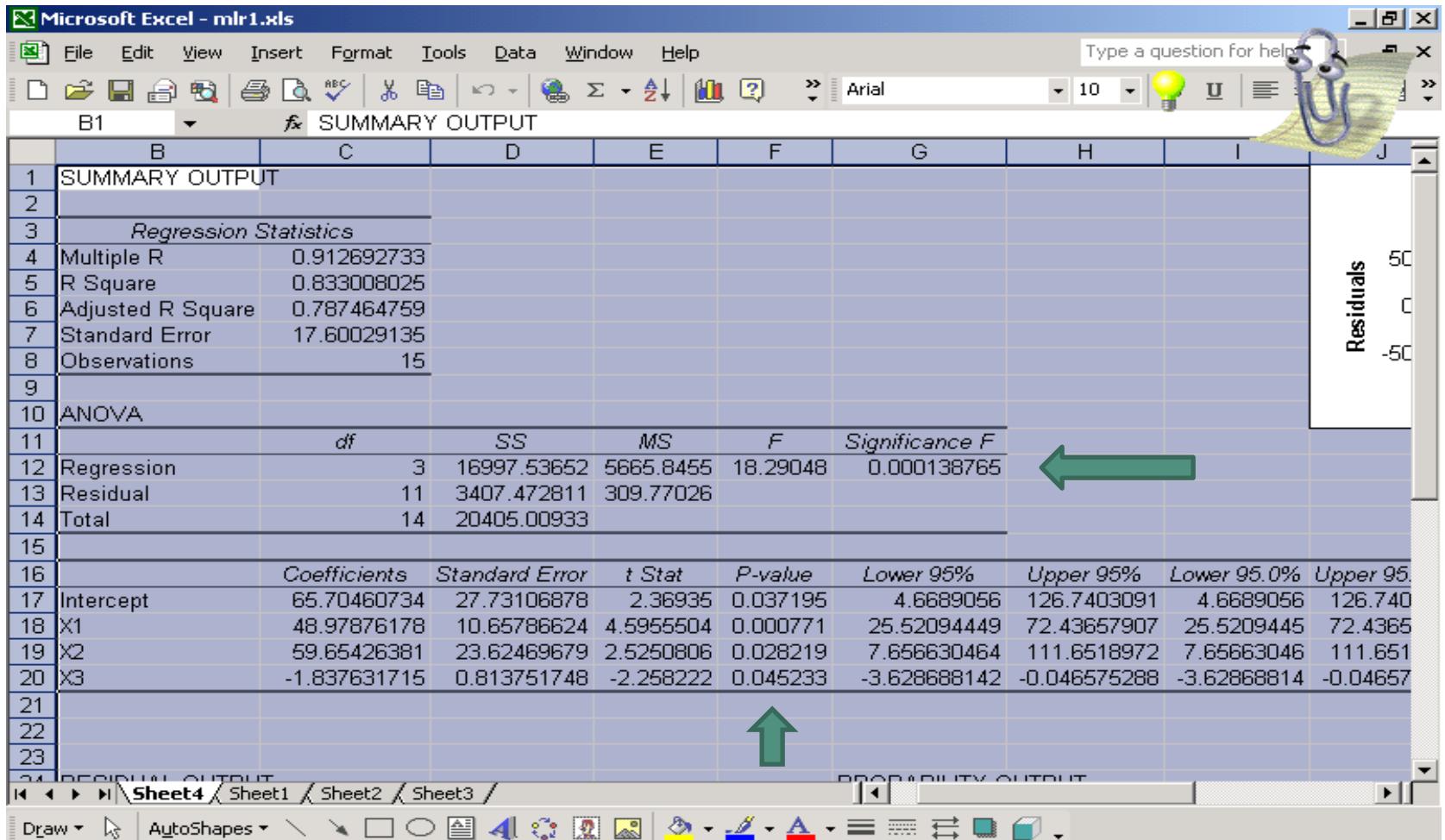
Note that if you include the labels (Y, X1, X2, X3) in the range, you must check the labels box.

Options

- Residuals is the difference between the observed value and the fitted value
- Standardized residuals are scaled to be invariant to the scale of the independent variable; for a good fit (no “outliers”) they should all be between 3 and -3
- Residual plots allow one to visually check the fit and see if they are random or have some pattern; if not random then the MLR model is not appropriate
- The Normal Probability Plot is used to check the normality assumption of the error term; there should be an even spread of data on either side of the 50% mark

Excel Output

[F and p-values are different for MLR]



Microsoft Excel - mlr1.xls

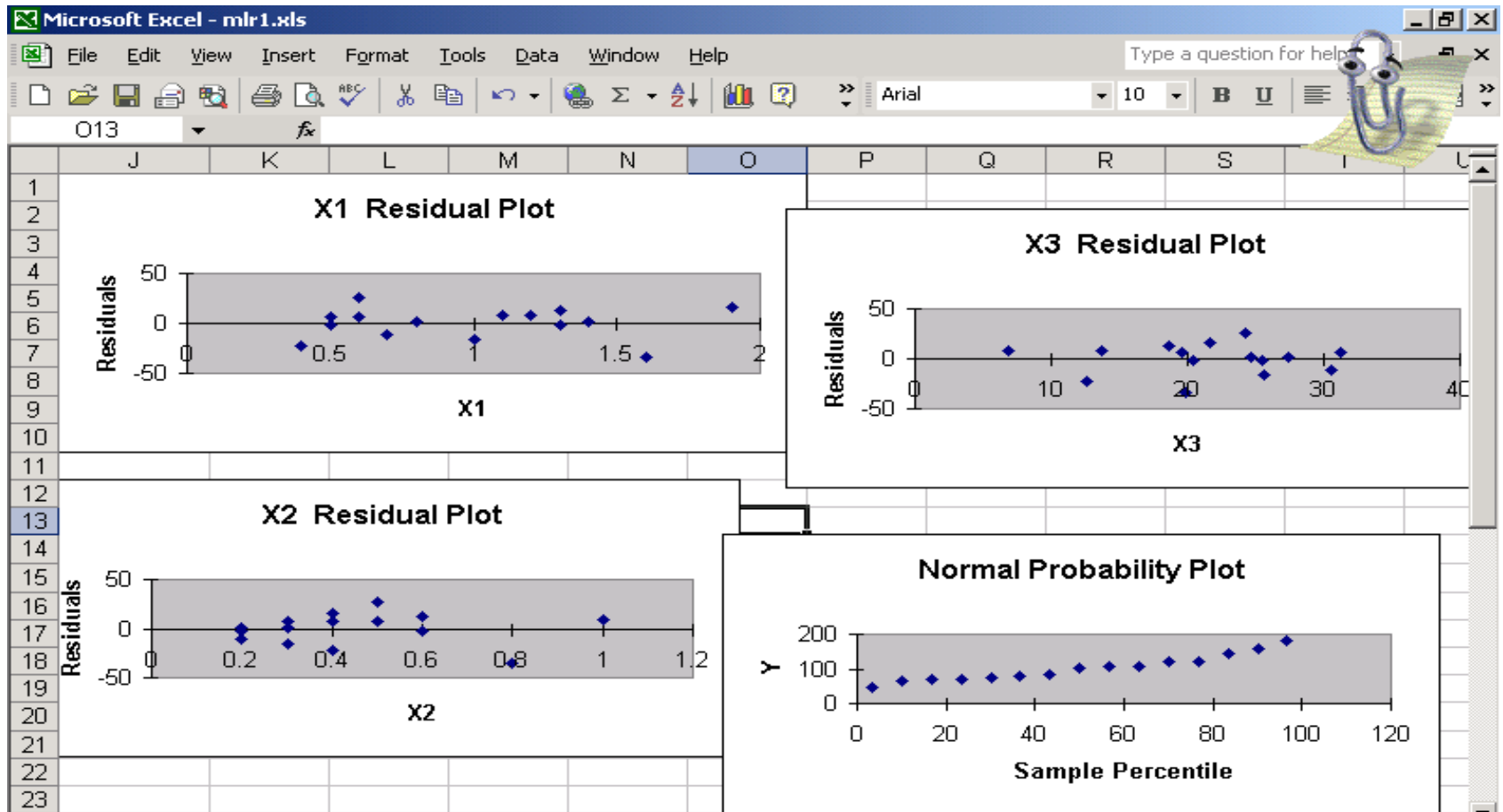
File Edit View Insert Format Tools Data Window Help

Summary OUTPUT

	B	C	D	E	F	G	H	I	J
1	SUMMARY OUTPUT								
2									
3	Regression Statistics								
4	Multiple R	0.912692733							
5	R Square	0.833008025							
6	Adjusted R Square	0.787464759							
7	Standard Error	17.60029135							
8	Observations	15							
9									
10	ANOVA								
11		df	SS	MS	F	Significance F			
12	Regression	3	16997.53652	5665.8455	18.29048	0.000138765			
13	Residual	11	3407.472811	309.77026					
14	Total	14	20405.00933						
15									
16		Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
17	Intercept	65.70460734	27.73106878	2.36935	0.037195	4.6689056	126.7403091	4.6689056	126.740
18	X1	48.97876178	10.65786624	4.5955504	0.000771	25.52094449	72.43657907	25.5209445	72.4365
19	X2	59.65426381	23.62469679	2.5250806	0.028219	7.656630464	111.6518972	7.65663046	111.651
20	X3	-1.837631715	0.813751748	-2.258222	0.045233	-3.628688142	-0.046575288	-3.62868814	-0.04657
21									
22									
23									
24	RESIDUAL OUTPUT								
25									

Sheet4 Sheet1 Sheet2 Sheet3

Residual Plots



Pizza Sales – What is the regression formula ?

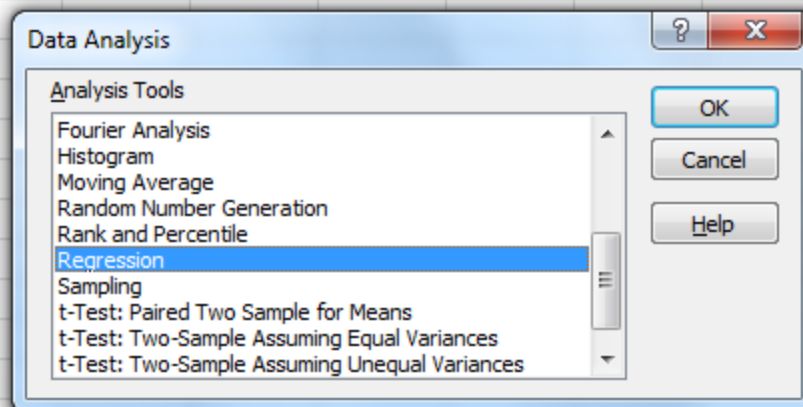
	A	B	C
1	PizzaStore	CampusSize	Sales
2	1	2	58
3	2	6	105
4	3	8	88
5	4	8	118
6	5	12	117
7	6	16	137
8	7	20	157
9	8	20	169
10	9	22	149
11	10	26	202

Wait....



Don't look ahead, until
you have your answer !

	A	B	C	D	E	F	G	H	I	J
1	PizzaStore	CampusSize	Sales							
2	1	2	58							
3	2	6	105							
4	3	8	88							
5	4	8	118							
6	5	12	117							
7	6	16	137							
8	7	20	157							
9	8	20	169							
10	9	22	149							
11	10	26	202							
12										
13										
14										
15										



	A	B	C	D	E	F	G	H	I	J	K
1	PizzaStore	CampusSize	Sales								
2	1	2	58								
3	2	6	105								
4	3	8	88								
5	4	8	118								
6	5	12	117								
7	6	16	137								
8	7	20	157								
9	8	20	169								
10	9	22	149								
11	10	26	202								
12											
13											
14											
15											
16											
17											
18											
19											
20											
21											
22											

Regression

Input

Input Y Range:

\$C\$1:\$C\$11

Input X Range:

\$B\$1:\$B\$11

☒ Labels

☐ Constant is Zero

☒ Confidence Level:

95

%

Output options

☐ Output Range:

☒ New Worksheet Ply:

☐ New Workbook

Residuals

☒ Residuals

☒ Residual Plots

☒ Standardized Residuals

☐ Line Fit Plots

Normal Probability

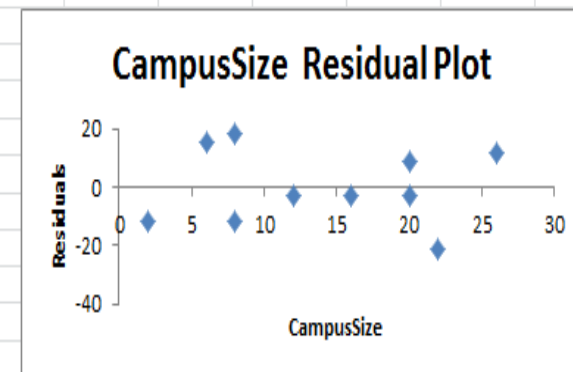
☒ Normal Probability Plots

OK

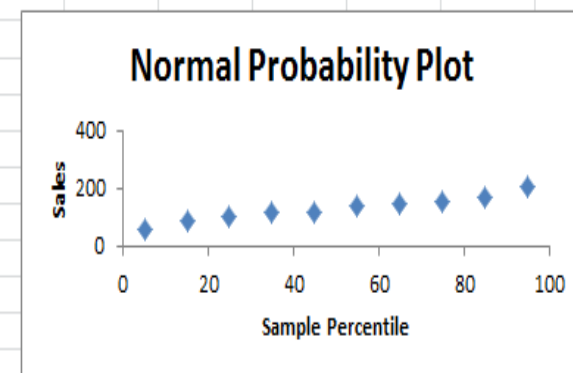
Cancel

Help





	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	SUMMARY OUTPUT															
2																
3	Regression Statistics															
4	Multiple R	0.950122955														
5	R Square	0.90273363														
6	Adjusted R Square	0.890575334														
7	Standard Error	13.82931669														
8	Observations	10														
9																
10	ANOVA															
11		df	SS	MS	F	Significance F										
12	Regression	1	14200	14200	74.24836601	2.54887E-05										
13	Residual	8	1530	191.25												
14	Total	9	15730													
15																
16		Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%							
17	Intercept	60	9.22603481	6.503336	0.000187444	38.72472558	81.27527	38.72473	81.27527							
18	CampusSize	5	0.580265238	8.616749	2.54887E-05	3.661905962	6.338094	3.661906	6.338094							
19																
20																
21																
22	RESIDUAL OUTPUT															
23																
24	Observation	Predicted Sales	Residuals	Standard Residuals												
25	1	70	-12	-0.92036												
26	2	90	15	1.150447												
27	3	100	-12	-0.92036												
28	4	100	18	1.380537												
29	5	120	-3	-0.23009												
30	6	140	-3	-0.23009												
31	7	160	-3	-0.23009												
32	8	160	9	0.690268												
33	9	170	-21	-1.61063												
34	10	190	12	0.920358												
35																
36																



$$Y(\text{sales}) = 60 + 5 \cdot X(\text{size})$$

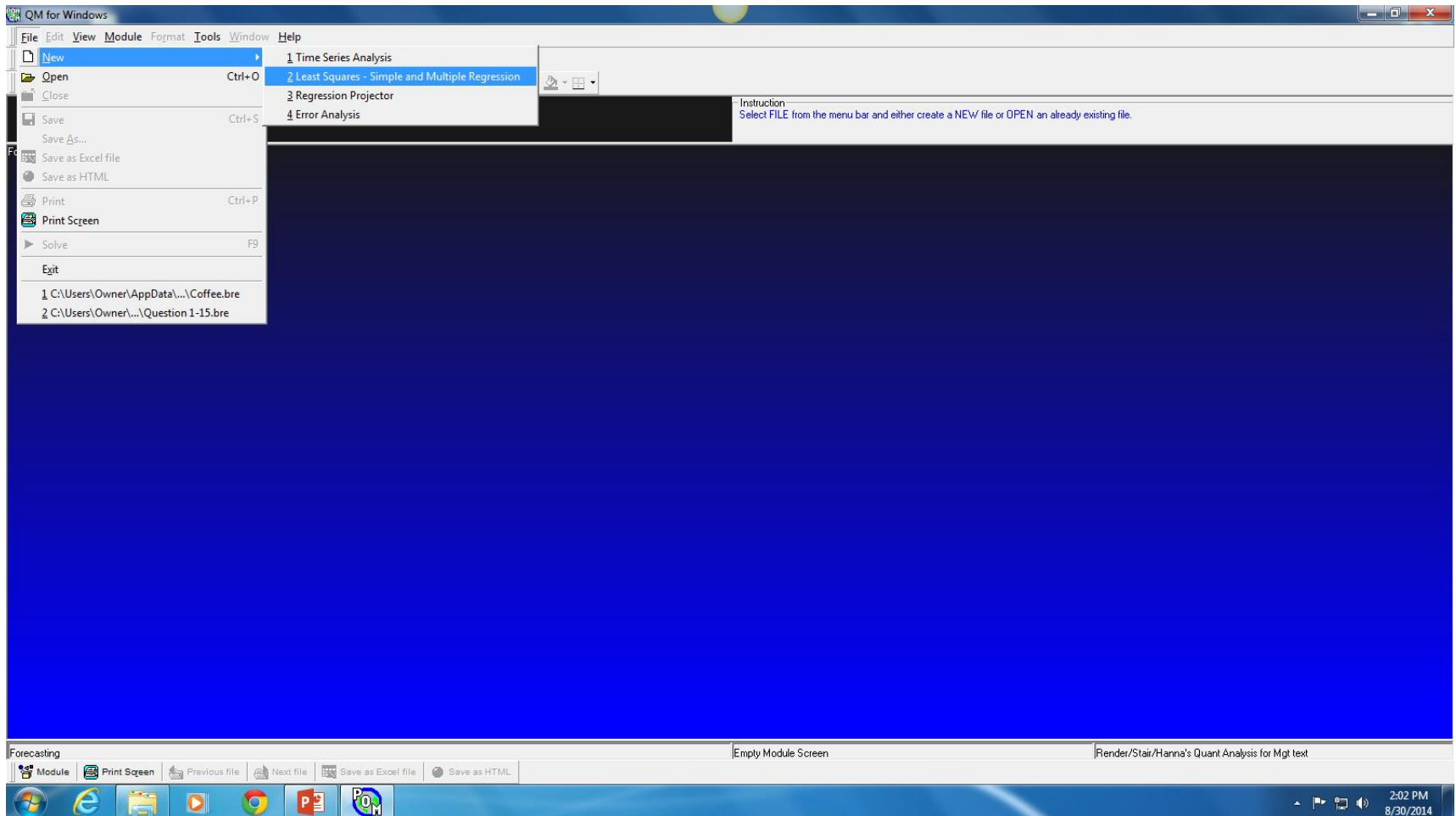


Excel Intercept and Slope Functions

C14		:	  	=INTERCEPT(C2:C11,B2:B11)				
	A	B	C	D	E	F		
1	PizzaStore	CampusSize	Sales					
2	1	2	58					
3	2	6	105					
4	3	8	88					
5	4	8	118					
6	5	12	117					
7	6	16	137					
8	7	20	157					
9	8	20	169					
10	9	22	149					
11	10	26	202					
12								
13								
14		Intercept	60					
15		Slope	5					
16								

Using QM

[need to go into “forecasting”]



Using QM (con't)

The screenshot shows a dialog box titled "Create data set for Forecasting/Least Squares - Simple and Multiple Regression". The dialog has a "Title" field containing "Pizza Sales" and a "Modify default title" button. Below this are two input fields: "Number of Observations" with a spinner set to 10, and "Number of Independent Variables" with a spinner set to 1. On the right, there are three tabs: "Row names", "Column names", and "Overview". The "Row names" tab is active, showing a list of radio button options: "Observation 1, Observation 2, Observation 3, ..." (selected), "a, b, c, d, e, ...", "A, B, C, D, E, ...", "1, 2, 3, 4, 5, ...", "January, February, March, April, ..." (with a sub-dropdown menu labeled "Click here to set start month"), and "Other". At the bottom are "Cancel", "Help", and "OK" buttons.

Create data set for Forecasting/Least Squares - Simple and Multiple Regression

Title: Pizza Sales Modify default title

Number of Observations: 10

Number of Independent Variables: 1

Row names | Column names | Overview

- ☒ Observation 1, Observation 2, Observation 3, ...
- ☐ a, b, c, d, e, ...
- ☐ A, B, C, D, E, ...
- ☐ 1, 2, 3, 4, 5, ...
- ☐ January, February, March, April, ...
Click here to set start month
- ☐ Other

Cancel Help OK

Using QM (con't)

QM for Windows - [Data Table]

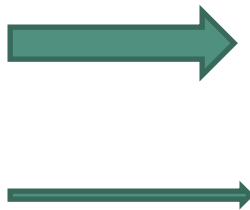
File Edit View Module Format Tools Window Help

Arial 8.25" B I U .000

Instruction
Enter the value of observation 10 for campusize. Any real value is permissible.

	Sales	CampusSize
Observation 1	58	2
Observation 2	105	6
Observation 3	88	8
Observation 4	118	8
Observation 5	117	12
Observation 6	137	16
Observation 7	157	20
Observation 8	169	20
Observation 9	149	22
Observation 10	202	26

Using QM (con't)



Forecasting Results	
Measure	Value
Error Measures	
Bias (Mean Error)	0
MAD (Mean Absolute Deviation)	10.8
MSE (Mean Squared Error)	153
Standard Error (denom=n-2=8)	13.829
MAPE (Mean Absolute Percent Error)	.096
Regression line	
Sales = 60	
+ 5 * CampusSize	
Statistics	
Correlation coefficient	.95
Coefficient of determination (r^2)	.903

Using QM (con't)

The screenshot displays the QM for Windows application. The main window is titled "Forecasting Results" and contains a table of statistical measures. A menu is open over the table, showing options: "Edit Data", "1 Forecasting Results", "2 Details and Error Analysis", "3 ANOVA Summary", "4 Sum of Squares Computations", and "5 Graph". The "ANOVA Summary" option is highlighted. The table lists various error measures and statistics for a regression model.

Measure	
Error Measures	
Bias (Mean Error)	0
MAD (Mean Absolute Deviation)	10.8
MSE (Mean Squared Error)	153
Standard Error (denom=n-2=8)	13.829
MAPE (Mean Absolute Percent Error)	.096
Regression line	
Sales = 60	
+ 5 * CampusSize	
Statistics	
Correlation coefficient	.95
Coefficient of determination (r^2)	.903

The bottom of the screen shows the Windows taskbar with the time 2:11 PM and date 8/30/2014. The application's status bar indicates the current screen is the "Solution Screen" for a "Forecasting/Least Squares - Simple and Multiple Regression" problem.

Using QM (con't)

QM for Windows

File Edit View Module Format Tools Window Help

100% Edit Data

Arial 8.25 B I U .000 Fix Dec 0.0

Instruction
There are more results available in additional windows. These may be opened by using the WINDOW option in the Main Menu.

Forecasting Results

Pizza Sales Summary

Measure	Value
Error Measures	
Bias (Mean Error)	0
MAD (Mean Absolute Deviation)	10.8
MSE (Mean Squared Error)	153
Standard Error (denom=n-2=8)	13.829
MAPE (Mean Absolute Percent Error)	.096
Regression line	
Sales = 60	
+ 5 * CampusSize	
Statistics	
Correlation coefficient	
Coefficient of determination (r^2)	

Output Table #3

Pizza Sales Solution

	Sum	Degrees of Freedom	Mean square
SSR (Sum of squares due to regression)	14200	1	14200
SSE (Sum of the squared error)	1530	8	191.25
SST (Sum of the squares total)	15730	9	
F statistic	74.248		
Probability	0		

Binary or Dummy Variables

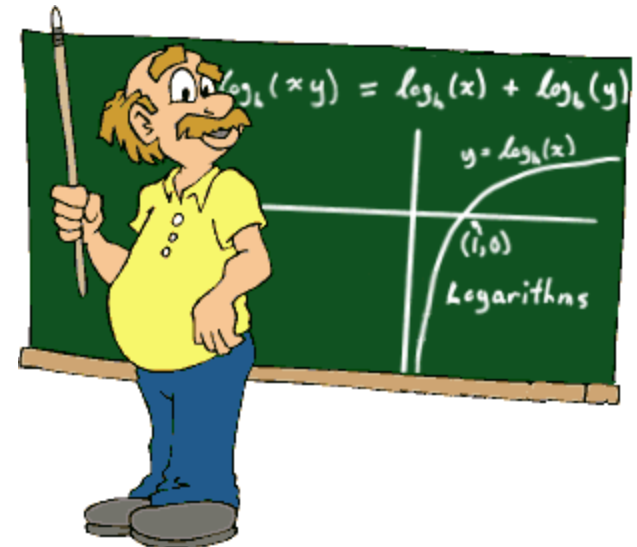
- *Binary* (or *dummy* or *indicator*) variables are special variables created for **qualitative data (nominal or ordinal)**
- A dummy variable is assigned a value of 1 if a particular condition is met and a value of 0 otherwise
- The number of dummy variables must equal one less than the number of categories of the qualitative variable

Realty Example

- In a model for house values, a regression model may have independent variables for square footage (X_1) and number of bedrooms (X_2)
- A better model could possibly be developed if information about the condition of the property was included:
$$X_3 = 1 \text{ if house is in excellent condition}$$
$$= 0 \text{ otherwise}$$
$$X_4 = 1 \text{ if house is in mint condition}$$
$$= 0 \text{ otherwise}$$
- Two dummy variables can be used to describe the three categories of condition (good, excellent, mint)
- No variable is needed for “good” condition since if both X_3 and $X_4 = 0$, the house must be in good condition

Transformation of Data

- Some non-linear data may be transformed into related linear data by applying a function to one or more independent variables:
 - Log function
 - Square root
 - Power function
 - Exponential

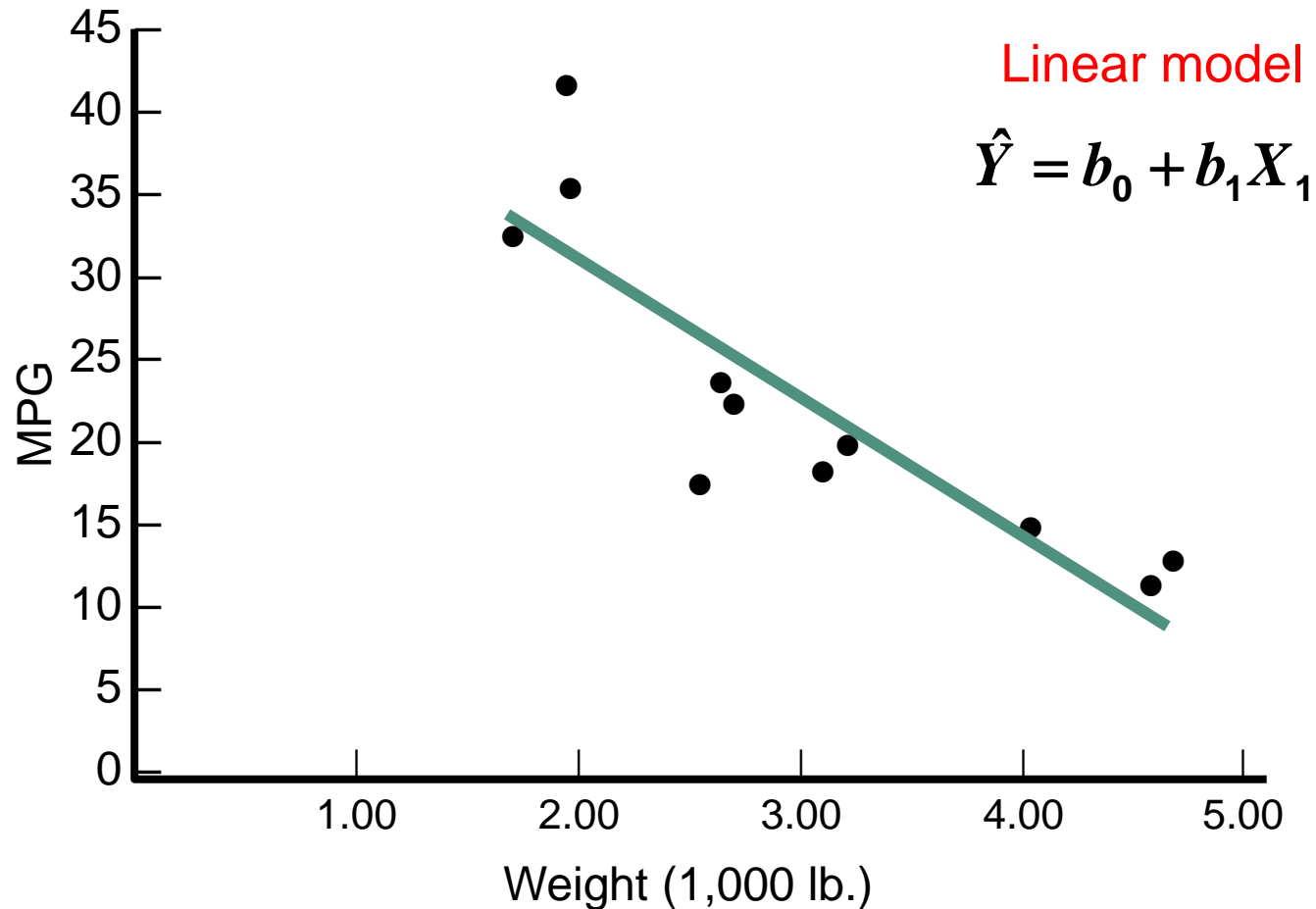


Transformation Example

- Engineers want to use regression analysis to improve fuel efficiency
- They have been asked to study the impact of weight on miles per gallon (MPG)

MPG	WEIGHT (1,000 LBS.)	MPG	WEIGHT (1,000 LBS.)
12	4.58	20	3.18
13	4.66	23	2.68
15	4.02	24	2.65
18	2.53	33	1.70
19	3.09	36	1.95
19	3.11	42	1.92

Transformation Example (con't)



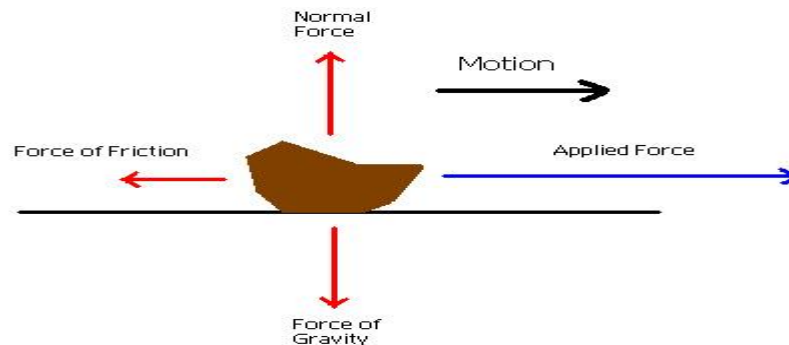
Transformation Example (con't)

	A	B +	C	D	E	F	G	H	I	J	K	L	M
1	Automobile Weight vs. MPG			SUMMARY OUTPUT									
2													
3	MPG (Y)	Weight (X1)		Regression Statistics									
4	12	4.58		Multiple R	0.8629								
5	13	4.66		R Square	0.7446								
6	15	4.02		Adjusted R Square	0.7190								
7	18	2.53		Standard Error	5.0076								
8	19	3.09		Observations	12								
9	19	3.11											
10	20	3.18		ANOVA									
11	23	2.68			df	SS	MS	F	Significance F				
12	24	2.65		Regression	1	730.9090	730.9090	29.1480	0.0003				
13	33	1.70		Residual	10	250.7577	25.0758						
14	36	1.95		Total	11	981.6667							
15	42	1.92											
16					Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%	
17				Intercept	47.6193	4.8132	9.8936	0.0000	36.8950	58.3437	36.8950	58.3437	
18				Weight	-8.2460	1.5273	-5.3989	0.0003	-11.6491	-4.8428	-11.6491	-4.8428	

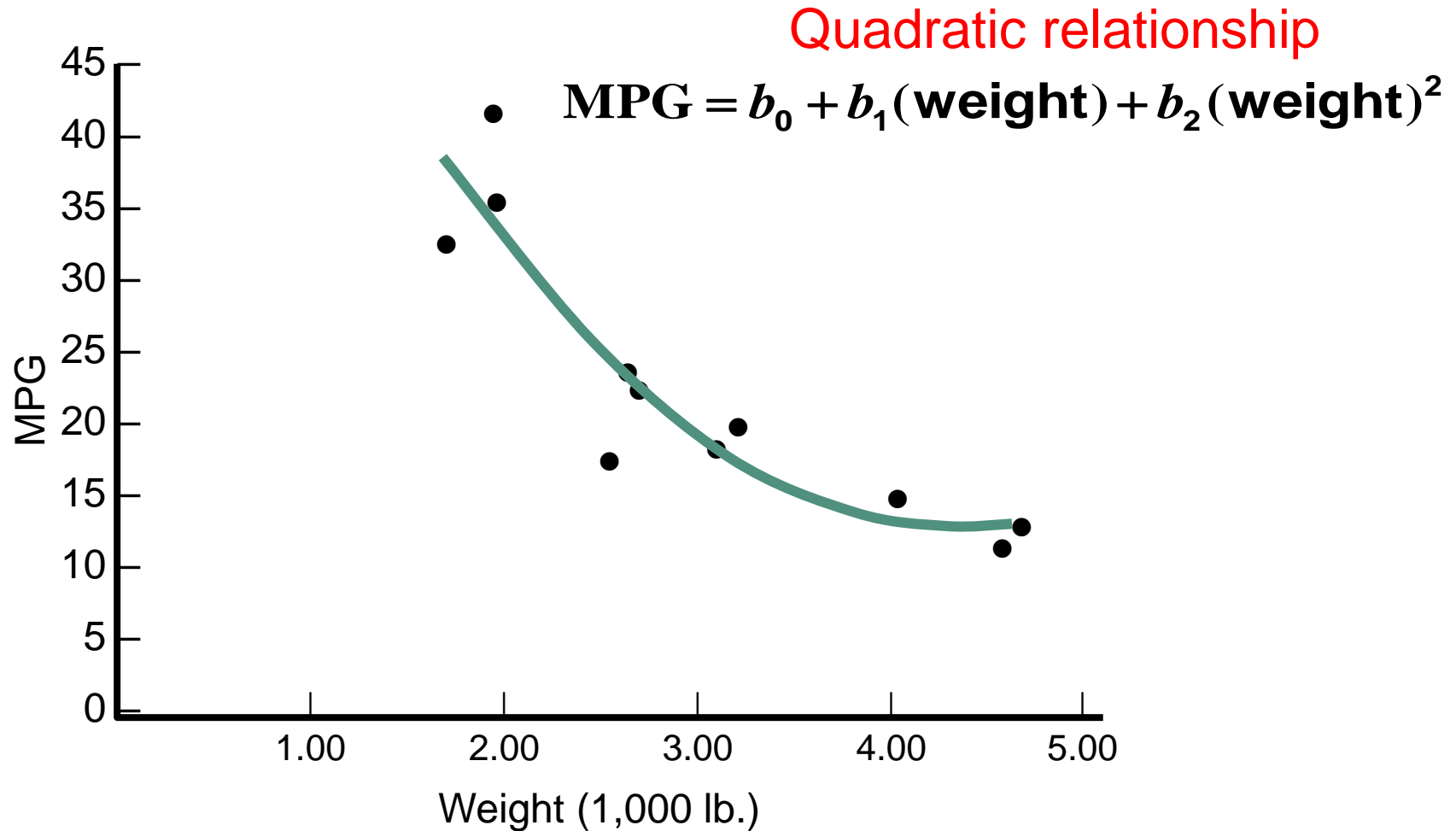
- A useful model with a small F -test for significance and a good r^2 value

Transformation Example (con't)

- What is a more likely relationship between MPG and weight ?
- Between the power needed to move an object of some weight ?
- Need force to accelerate an object ($F=MA$) and need force to overcome friction ($F=Mg$)



Transformation Example (con't)



Transformation Example (con't)

- The nonlinear model is a quadratic model
- The easiest way to work with this model is to **develop a new variable**

$$X_2 = (\text{weight})^2$$

- This gives us a model that can be solved with linear regression software

$$\hat{Y} = b_0 + b_1X_1 + b_2X_2$$

Transformation Example (con't)

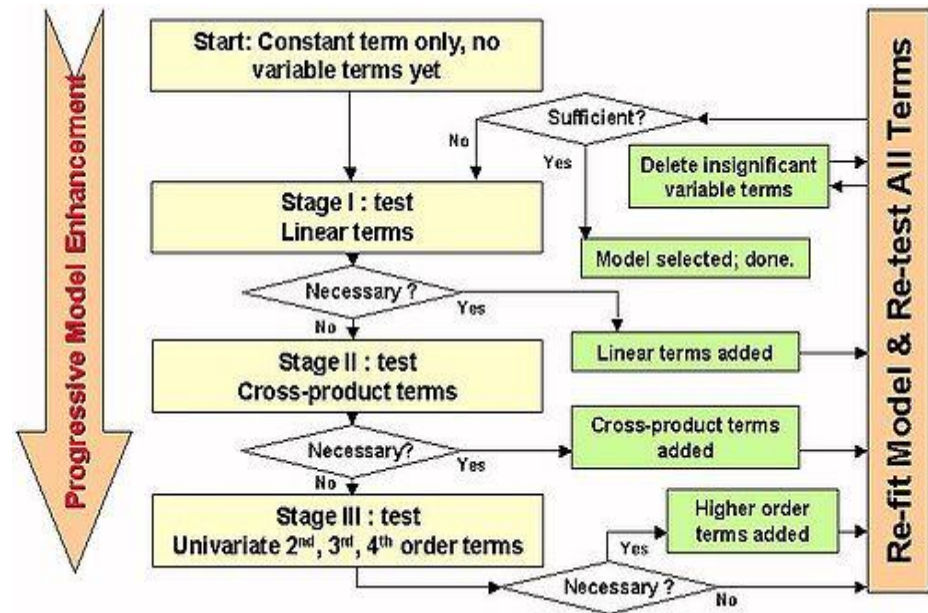
	A	B	C	D	E	F	G	H	I	J	K	L
1	Automobile Weight vs. MPG			SUMMARY OUTPUT								
2												
3	MPG (Y)	Weight (X1)	WeightSq. (X2)	Regression Statistics								
4	12	4.58	20.98	Multiple R	0.9208							
5	13	4.66	21.72	R Square	0.8478							
6	15	4.02	16.16	Adjusted R Square	0.8140							
7	18	2.53	6.40	Standard Error	4.0745							
8	19	3.09	9.55	Observations	12							
9	19	3.11	9.67									
10	20	3.18	10.11	ANOVA								
11	23	2.68	7.18		df	SS	MS	F	Significance F			
12	24	2.65	7.02	Regression	2	832.2557	416.1278	25.0661	0.0002			
13	33	1.70	2.89	Residual	9	149.4110	16.6012					
14	36	1.95	3.80	Total	11	981.6667						
15	42	1.92	3.69									
16					Coefficient	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
17				Intercept	79.7888	13.5962	5.8685	0.0002	49.0321	110.5454	49.0321	110.5454
18				Weight	-30.2224	8.9809	-3.3652	0.0083	-50.5386	-9.9061	-50.5386	-9.9061
19				Weight2	3.4124	1.3811	2.4708	0.0355	0.2881	6.5367	0.2881	6.5367

$$\hat{Y} = 79.8 - 30.2X_1 + 3.4X_2$$

- A better model with a smaller F -test for significance and a larger adjusted r^2 value

Stepwise MLR

- Brings variables into the regression one at a time to determine which variables have the most impact and which variables are really necessary
- Uses adjusted R squared to judge “improvement”
- Can also investigate variable relationships



Regression Pitfalls

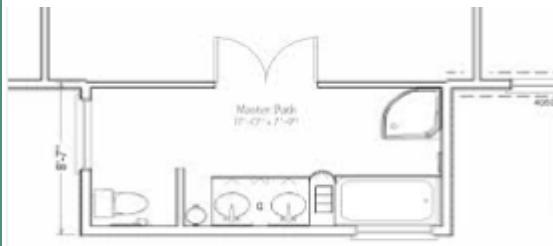
- If the **assumptions are not met**, the statistical test may not be valid
- **Correlation does not necessarily mean causation**
- **Multicollinearity** makes interpreting coefficients problematic, but the model may still be good
- Using a regression model **beyond the range of X** is questionable, the relationship may not hold outside the sample data
- A linear relationship may not be the best relationship, even if the F -test returns an acceptable value
- A **nonlinear** relationship can exist even if a linear relationship does not
- Just because a relationship is statistically significant doesn't mean it has any practical value

References

- [Data Analysis Using Regression and Multilevel/Hierarchical Models](#) by [Andrew Gelman](#) and [Jennifer Hill](#)
- [Applied Regression Analysis and Generalized Linear Models](#) by [John Fox](#)
- [Regression Analysis by Example \(Wiley Series in Probability and Statistics\)](#) by [Samprit Chatterjee](#) and Ali S. Hadi

Homework

- Textbook Chapter 4
- Quiz on these slides and Chapter 4
- Questions to be answered: 1, 2, 3, 5 from Chapter 4
- Project Two →



Project 2



- Betty Byte lives in a 4br/2.5ba 2900 sq ft house on .6 acres
- She has obtained recent sales data on comparable houses in her subdivision (shown on the next slide)
- She is considering adding another bath to her house which would be 300 sq ft
- What is the most she should spend on adding that extra bath ?
 - Use MLR to find the relevant value formula and how much adding one bath at 300 sq ft to her house would increase the value of her house

Project (con't)



Value	Acres	Sq Ft	BRs	Baths
253000	0.5	3000	3	2
310000	0.6	3400	4	3.5
260000	0.4	3100	3	2.5
340000	0.7	3600	5	2.5
320000	0.8	3200	4	3
305000	0.55	3300	5	2.5
285000	0.65	2900	4	2.5
278000	0.45	3200	4	2.5
325000	0.75	3300	4	3
315000	0.5	3600	5	3

Regression Derivation

[Khan Academy Videos for Detailed Derivation]

