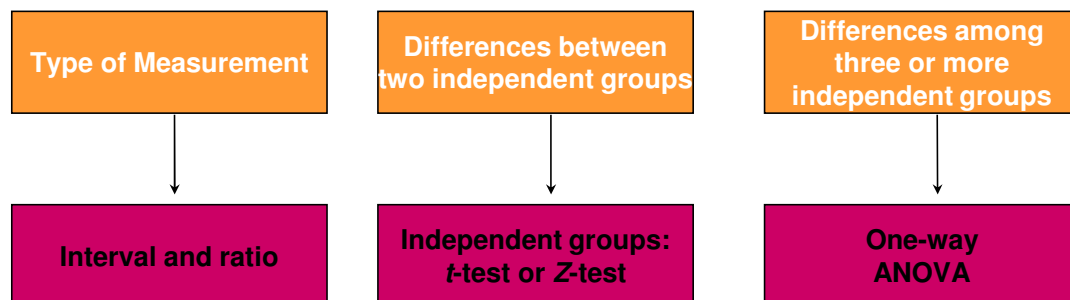


Research Methods

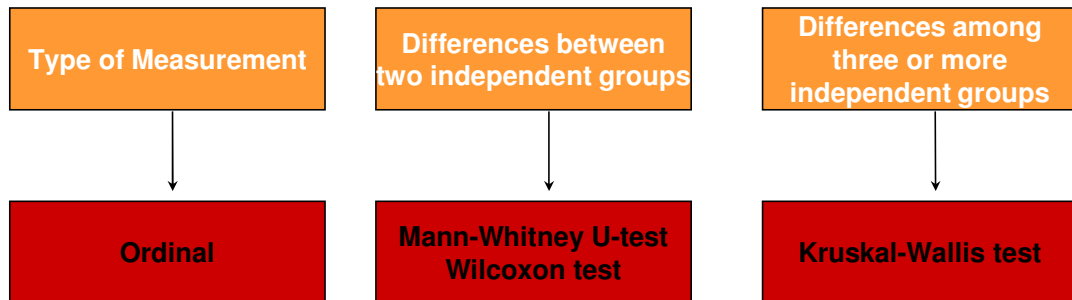
William G. Zikmund

Bivariate Analysis - Tests of Differences

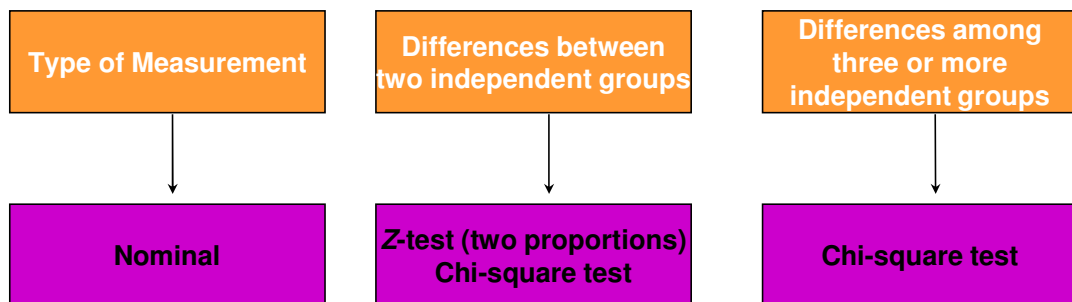
Common Bivariate Tests

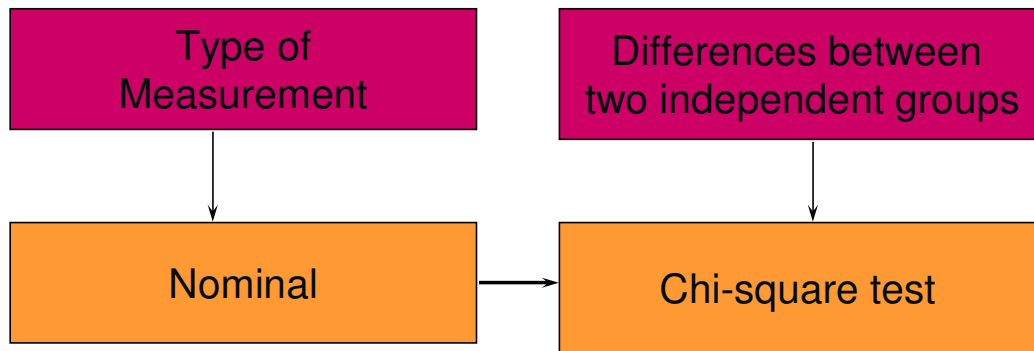


Common Bivariate Tests



Common Bivariate Tests





Differences Between Groups

- Contingency Tables
- Cross-Tabulation
- Chi-Square allows testing for significant differences between groups
- “Goodness of Fit”

Chi-Square Test

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

χ^2 = chi-square statistics

O_i = observed frequency in the i^{th} cell

E_i = expected frequency on the i^{th} cell

Chi-Square Test

$$E_{ij} = \frac{R_i C_j}{n}$$

R_i = total observed frequency in the i^{th} row

C_j = total observed frequency in the j^{th} column

n = sample size

Degrees of Freedom

$$(R-1)(C-1)=(2-1)(2-1)=1$$

Health Economics Research Method 200

Degrees of Freedom

$$\text{d.f.}=(R-1)(C-1)$$

Awareness of Tire Manufacturer's Brand

	Men	Women	Total
Aware	50	10	60
Unaware	$\frac{15}{65}$	$\frac{25}{35}$	$\frac{40}{100}$

Chi-Square Test: Differences Among Groups Example

$$X^2 = \frac{(50 - 39)^2}{39} + \frac{(10 - 21)^2}{21} + \frac{(15 - 26)^2}{26} + \frac{(25 - 14)^2}{14}$$

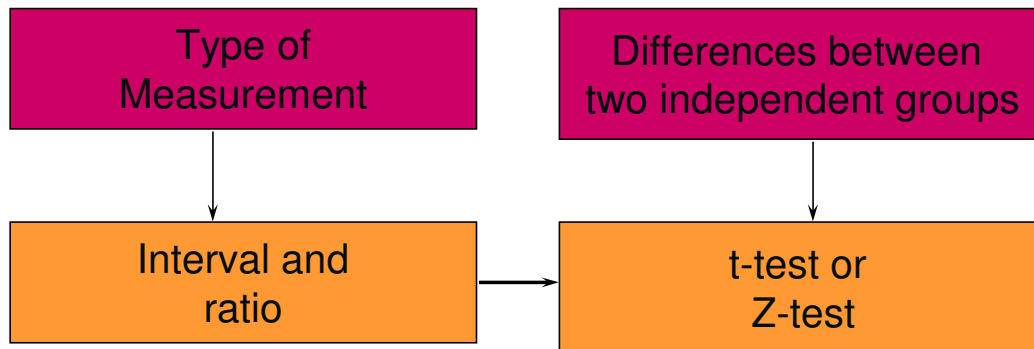
$$\chi^2 = 3.102 + 5.762 + 4.654 + 8.643 =$$

$$\chi^2 = 22.161$$

$$d.f. = (R - 1)(C - 1)$$

$$d.f. = (2 - 1)(2 - 1) = 1$$

$\chi^2 = 3.84$ with 1 d.f.



Differences Between Groups when Comparing Means

- Ratio scaled dependent variables
- t-test
 - When groups are small
 - When population standard deviation is unknown
- z-test
 - When groups are large

Null Hypothesis About Mean Differences Between Groups

$$\mu_1 - \mu_2$$

OR

$$\mu_1 - \mu_2 = 0$$

t-Test for Difference of Means

$$t = \frac{\text{mean 1} - \text{mean 2}}{\text{Variability of random means}}$$

t-Test for Difference of Means

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{\bar{X}_1 - \bar{X}_2}}$$

\bar{X}_1 = mean for Group 1

\bar{X}_2 = mean for Group 2

$S_{\bar{X}_1 - \bar{X}_2}$ = the pooled or combined standard error of difference between means.

t-Test for Difference of Means

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{\bar{X}_1 - \bar{X}_2}}$$

t-Test for Difference of Means

\bar{X}_1 = mean for Group 1

\bar{X}_2 = mean for Group 2

$S_{\bar{X}_1 - \bar{X}_2}$ = the pooled or combined standard error of difference between means.

Pooled Estimate of the Standard Error

$$S_{\bar{X}_1 - \bar{X}_2} = \sqrt{\left(\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2} \right) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

Pooled Estimate of the Standard Error

S_1^2 = the variance of Group 1

S_2^2 = the variance of Group 2

n_1 = the sample size of Group 1

n_2 = the sample size of Group 2

Pooled Estimate of the Standard Error

t-test for the Difference of Means

$$S_{\bar{X}_1 - \bar{X}_2} = \sqrt{\left(\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2} \right) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

S_1^2 = the variance of Group 1

S_2^2 = the variance of Group 2

n_1 = the sample size of Group 1

n_2 = the sample size of Group 2

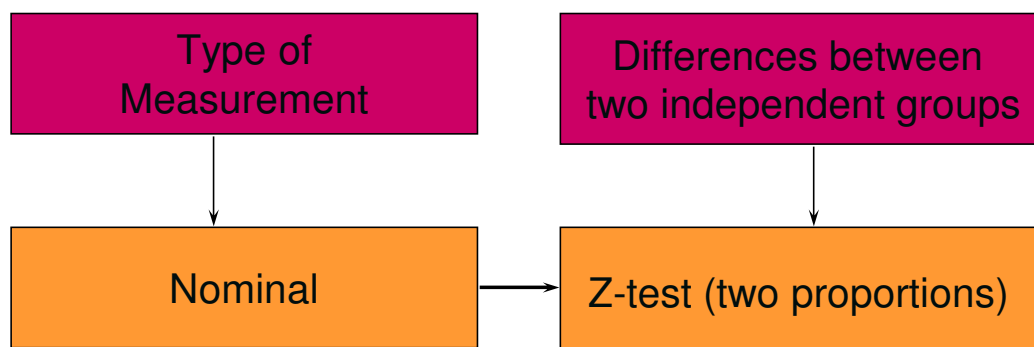
Degrees of Freedom

- d.f. = $n - k$
- where:
 - $n = n_1 + n_2$
 - $k =$ number of groups

t-Test for Difference of Means Example

$$s_{\bar{x}_1 - \bar{x}_2} = \sqrt{\left(\frac{(20)(2.1)^2 + (13)(2.6)^2}{33} \right) \left(\frac{1}{21} + \frac{1}{14} \right)}$$
$$= .797$$

$$t = \frac{16.5 - 12.2}{.797} = \frac{4.3}{.797}$$
$$= 5.395$$



Comparing Two Groups when Comparing Proportions

- Percentage Comparisons
- Sample Proportion - P
- Population Proportion - Π

Differences Between Two Groups when Comparing Proportions

The hypothesis is:

$$H_0: \Pi_1 = \Pi_2$$

may be restated as:

$$H_0: \Pi_1 - \Pi_2 = 0$$

Z-Test for Differences of Proportions

$$H_o : \pi_1 = \pi_2$$

or

$$H_o : \pi_1 - \pi_2 = 0$$

Z-Test for Differences of Proportions

$$Z = \frac{(p_1 - p_2) - (\pi_1 - \pi_2)}{S_{p_1 - p_2}}$$

Z-Test for Differences of Proportions

p_1 = sample portion of successes in Group 1

p_2 = sample portion of successes in Group 2

$(\pi_1 - \pi_2)$ = hypothesized population proportion 1
minus hypothesized population
proportion 2

$S_{p_1-p_2}$ = pooled estimate of the standard errors of
difference of proportions

Z-Test for Differences of Proportions

$$S_{p_1-p_2} = \sqrt{\bar{p}\bar{q} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

Z-Test for Differences of Proportions

\bar{p} = pooled estimate of proportion of success in a sample of both groups

$\bar{q} = (1 - \bar{p})$ or a pooled estimate of proportion of failures in a sample of both groups

n_1 = sample size for group 1

n_2 = sample size for group 2

Z-Test for Differences of Proportions

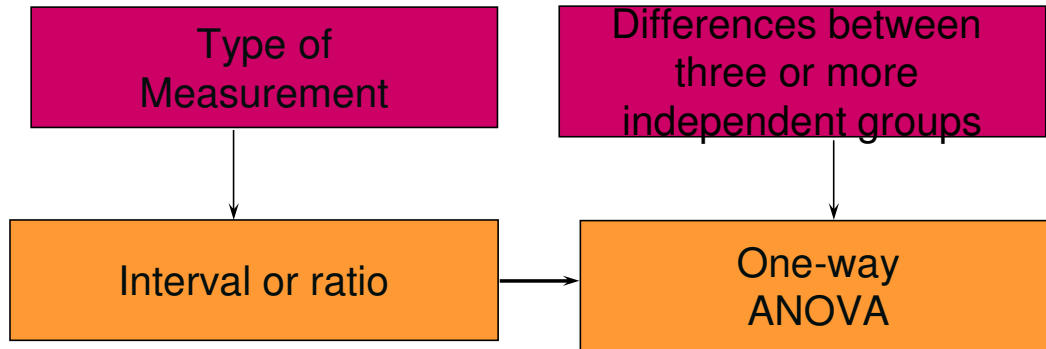
$$\bar{p} = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}$$

Z-Test for Differences of Proportions

$$S_{p_1-p_2} = \sqrt{(.375)(.625)\left(\frac{1}{100} + \frac{1}{100}\right)}$$
$$= .068$$

A Z-Test for Differences of Proportions

$$\bar{p} = \frac{(100)(.35) + (100)(.4)}{100 + 100}$$
$$= .375$$



Analysis of Variance

Hypothesis when comparing three groups

$$\mu_1 = \mu_2 = \mu_3$$

Analysis of Variance F-Ratio

$$F = \frac{\text{Variance – between – groups}}{\text{Variance – within – groups}}$$

Analysis of Variance Sum of Squares

$$SS_{\text{total}} = SS_{\text{within}} + SS_{\text{between}}$$

Analysis of Variance Sum of Squares Total

$$SS_{\text{total}} = \sum_{i=1}^n \sum_{j=1}^c (X_{ij} - \bar{\bar{X}})^2$$

Analysis of Variance Sum of Squares

X_{ij} = individual scores, i.e., the i^{th} observation or test unit in the j^{th} group

$\bar{\bar{X}}$ = grand mean

n = number of all observations or test units in a group

c = number of j^{th} groups (or columns)

Analysis of Variance Sum of Squares Within

$$SS_{\text{within}} = \sum_{i=1}^n \sum_{j=1}^c (X_{ij} - \bar{X}_j)^2$$

Analysis of Variance Sum of Squares Within

X_{ij} = individual scores, i.e., the i^{th} observation or test unit in the j^{th} group

\bar{X} = grand mean

n = number of all observations or test units in a group

c = number of j^{th} groups (or columns)

Analysis of Variance Sum of Squares Between

$$SS_{\text{between}} = \sum_{j=1}^n n_j (\bar{X}_j - \bar{\bar{X}})^2$$

Analysis of Variance Sum of squares Between

$X_j =$ individual scores, i.e., the i^{th} observation or test unit in the j^{th} group

$\bar{\bar{X}} =$ grand mean

$n_j =$ number of all observations or test units in a group

Analysis of Variance Mean Squares Between

$$MS_{between} = \frac{SS_{between}}{c - 1}$$

Analysis of Variance Mean Square Within

$$MS_{within} = \frac{SS_{within}}{cn - c}$$

Analysis of Variance F-Ratio

$$F = \frac{MS_{between}}{MS_{within}}$$

A Test Market Experiment on Pricing

Sales in Units (thousands)

	Regular Price \$.99	Reduced Price \$.89	Cents-Off Coupon Regular Price
Test Market A, B, or C	130	145	153
Test Market D, E, or F	118	143	129
Test Market G, H, or I	87	120	96
Test Market J, K, or L	84	131	99
Mean	$X_1=104.75$	$X_2=134.75$	$X_1=119.25$
Grand Mean	$X=119.58$		

ANOVA Summary Table Source of Variation

- Between groups
- Sum of squares
 - SS_{between}
- Degrees of freedom
 - $c-1$ where c =number of groups
- Mean squared- MS_{between}
 - $SS_{\text{between}}/c-1$

ANOVA Summary Table Source of Variation

- Within groups
- Sum of squares
 - SS_{within}
- Degrees of freedom
 - $cn-c$ where c =number of groups, n = number of observations in a group
- Mean squared- MS_{within}
 - $SS_{\text{within}}/cn-c$

ANOVA Summary Table

Source of Variation

- Total
- Sum of Squares
 - SStotal
- Degrees of Freedom
 - $cn-1$ where c =number of groups, n = number of observations in a group

$$F = \frac{MS_{BETWEEN}}{MS_{WITHIN}}$$

The screenshot shows a Microsoft Excel spreadsheet titled "MLB t-test". The data is organized into two columns: "Beer" (Column A) and "Soda" (Column B). The values for Beer are: 12, 16, 20, 20, 16, 16, 16, 24, 20, 16, 24, 14, 14, 14, 16, 16, 21, 16, 22, 20. The values for Soda are: 14, 18, 20, 16, 16, 12, 14, 14, 16, 16, 16, 14, 14, 14, 16, 16, 20, 12, 16, 16, 24. A t-test summary table is displayed in columns D through F, with the following data:

	Variable 1	Variable 2
Mean	16.83333333	15.8
Variance	10.35057471	7.820689655
Observations	30	30
Pearson Correlation	0.264451682	
Hypothesized Mean Difference	0	
df	29	
t Stat	1.545410551	
P(T<=t) one-tail	0.06654553	
t Critical one-tail	1.699127097	
P(T<=t) two-tail	0.13309106	
t Critical two-tail	2.045230758	

Research Methods

William G. Zikmund

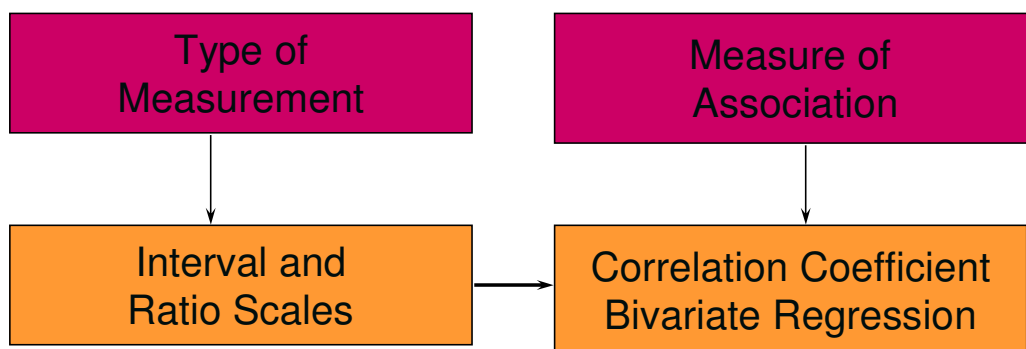
Bivariate Analysis: Measures of Associations

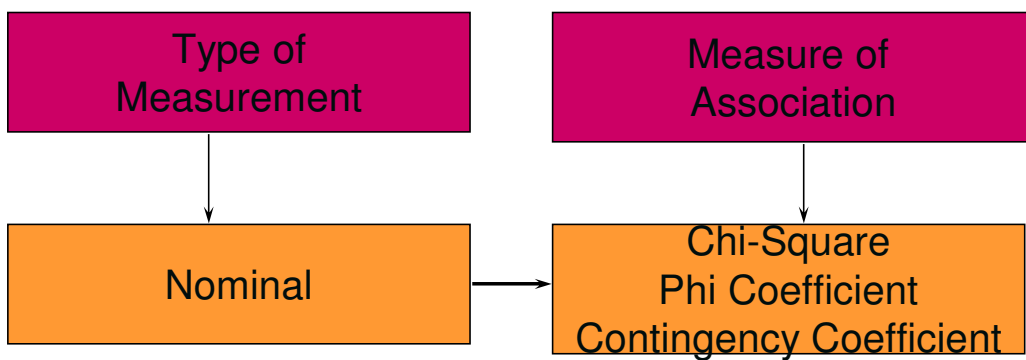
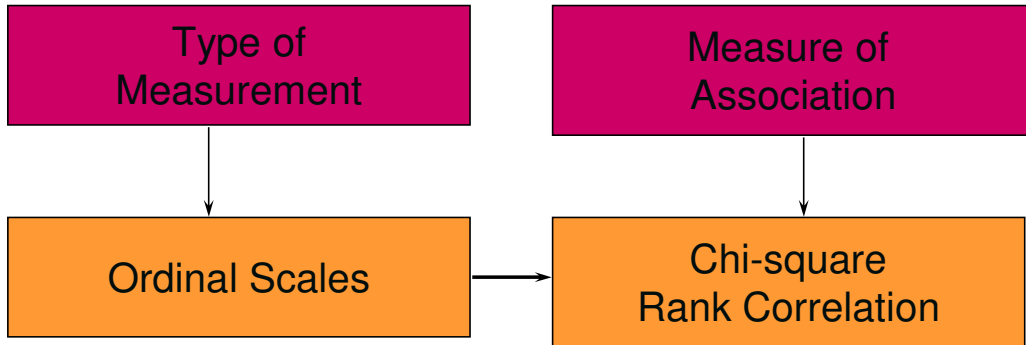
Measures of Association

- A general term that refers to a number of bivariate statistical techniques used to measure the strength of a relationship between two variables.

Relationships Among Variables

- Correlation analysis
- Bivariate regression analysis





Correlation Coefficient

- A statistical measure of the covariation or association between two variables.
- Are dollar sales associated with advertising dollar expenditures?

The Correlation coefficient for two variables, X and Y is

$$r_{xy}$$

Correlation Coefficient

- r
- r ranges from +1 to -1
- $r = +1$ a perfect positive linear relationship
- $r = -1$ a perfect negative linear relationship
- $r = 0$ indicates no correlation

Simple Correlation Coefficient

$$r_{xy} = r_{yx} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}}$$

Simple Correlation Coefficient

$$r_{xy} = r_{yx} = \frac{\sigma_{xy}}{\sqrt{\sigma_x^2 \sigma_y^2}}$$

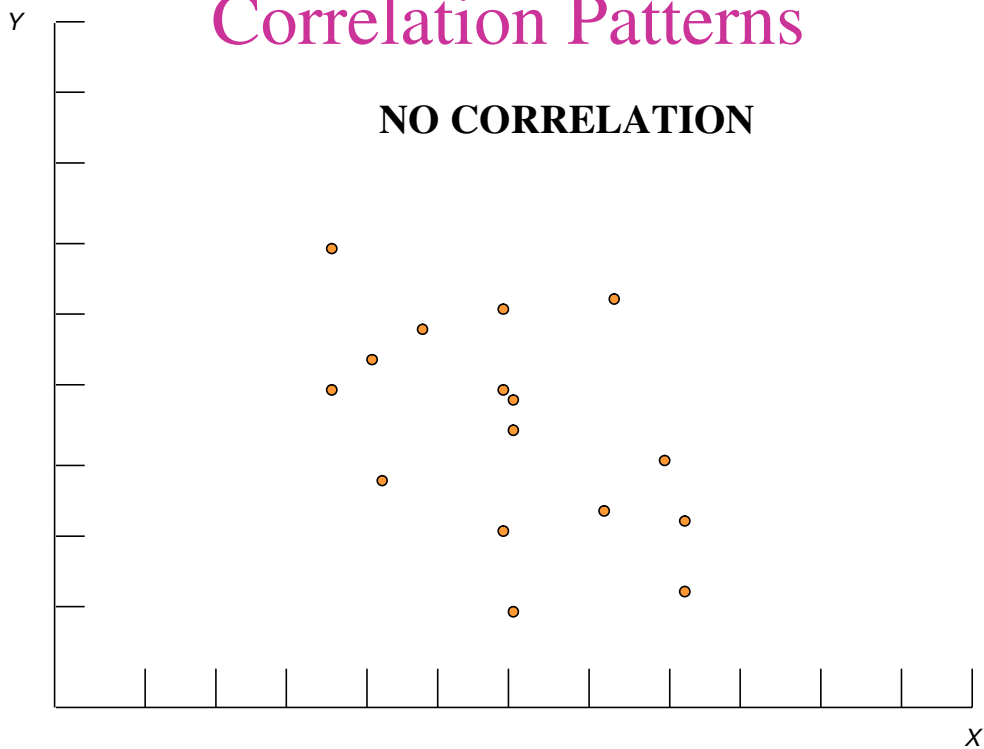
Simple Correlation Coefficient Alternative Method

σ_x^2 = Variance of X

σ_y^2 = Variance of Y

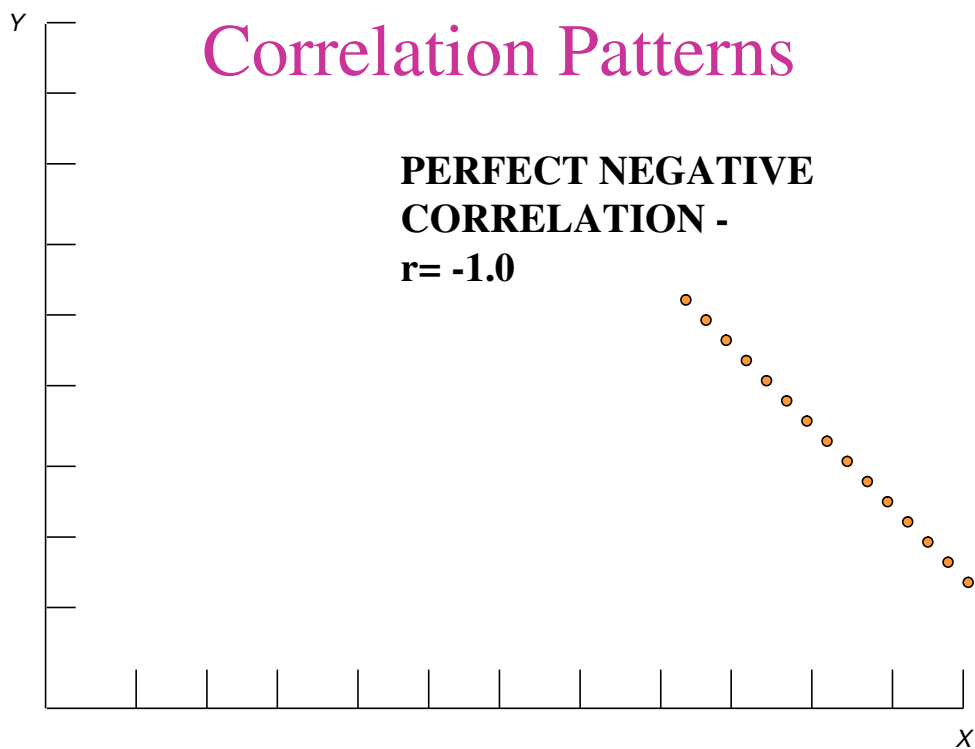
σ_{xy} = Covariance of X and Y

Correlation Patterns



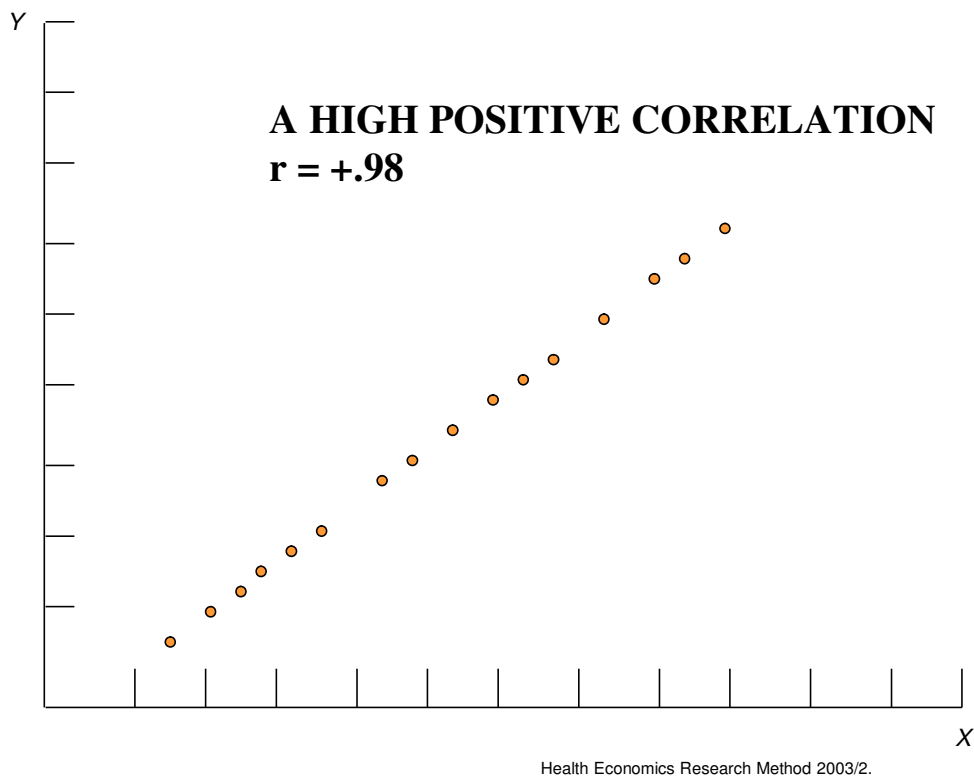
Health Economics Research Method 2003/2.

Correlation Patterns



Health Economics Research Method 2003/2.

Correlation Patterns



Calculation of r

$$\begin{aligned} r &= \frac{-6.3389}{\sqrt{(17.837)(5.589)}} \\ &= \frac{-6.3389}{\sqrt{99.712}} = -.635 \end{aligned}$$

Coefficient of Determination

$$r^2 = \frac{\textit{Explained variance}}{\textit{Total Variance}}$$

Correlation Does Not Mean Causation



- High correlation
- Rooster's crow and the rising of the sun
 - Rooster does not cause the sun to rise.
- Teachers' salaries and the consumption of liquor
 - Covary because they are both influenced by a third variable

Correlation Matrix

- The standard form for reporting correlational results.

Correlation Matrix

	Var1	Var2	Var3
Var1	1.0	0.45	0.31
Var2	0.45	1.0	0.10
Var3	0.31	0.10	1.0

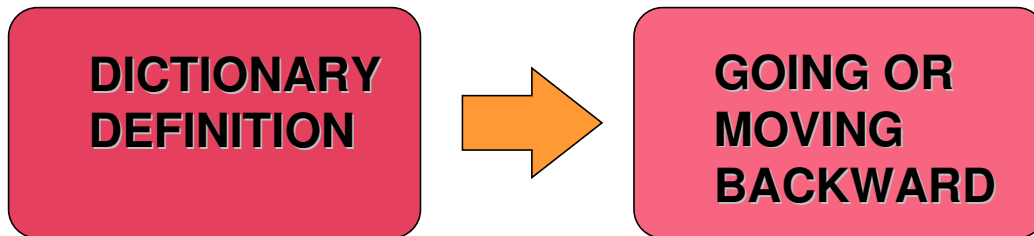
Walkup's First Laws of Statistics

- Law No. 1
 - Everything correlates with everything, especially when the same individual defines the variables to be correlated.
- Law No. 2
 - It won't help very much to find a good correlation between the variable you are interested in and some other variable that you don't understand any better.

Walkup's First Laws of Statistics

- Law No. 3
 - Unless you can think of a logical reason why two variables should be connected as cause and effect, it doesn't help much to find a correlation between them. In Columbus, Ohio, the mean monthly rainfall correlates very nicely with the number of letters in the names of the months!

Regression



Going back to previous conditions

- Tall men's sons

Bivariate Regression

- A measure of linear association that investigates a straight line relationship
- Useful in forecasting

Bivariate Linear Regression

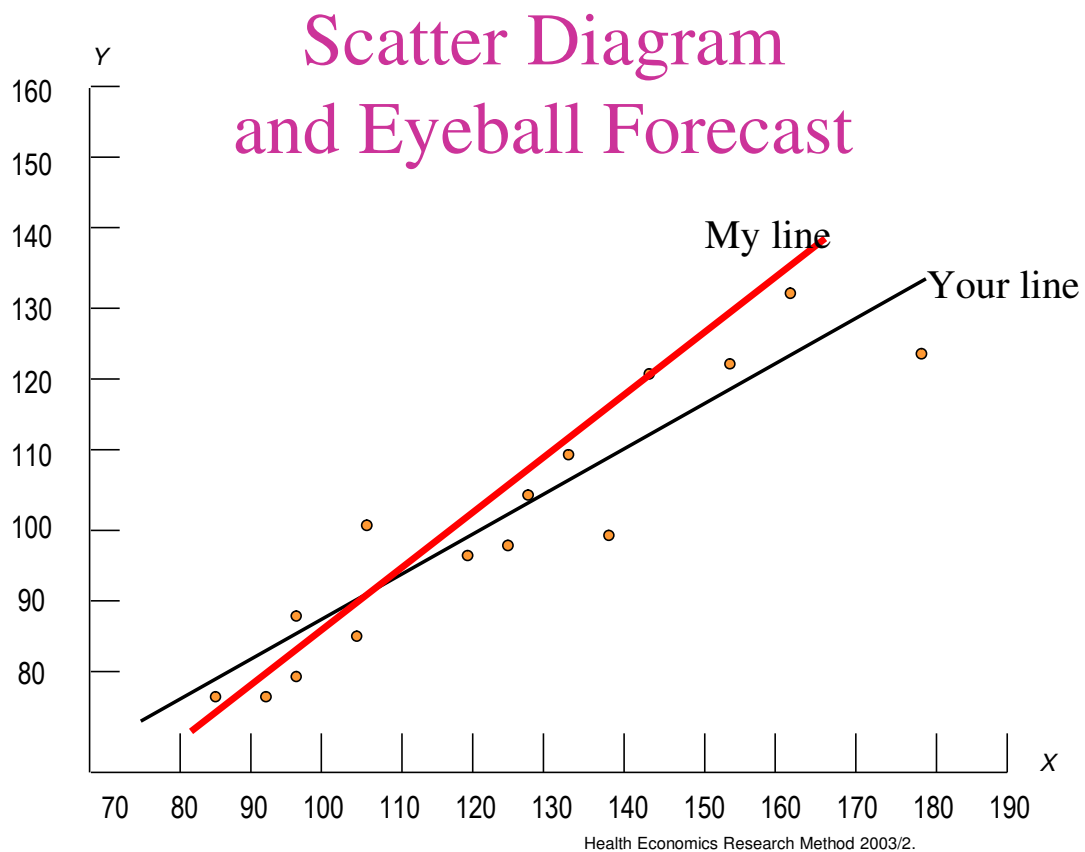
- A measure of linear association that investigates a straight-line relationship
- $Y = a + bX$
- where
- Y is the dependent variable
- X is the independent variable
- a and b are two constants to be estimated

Y intercept

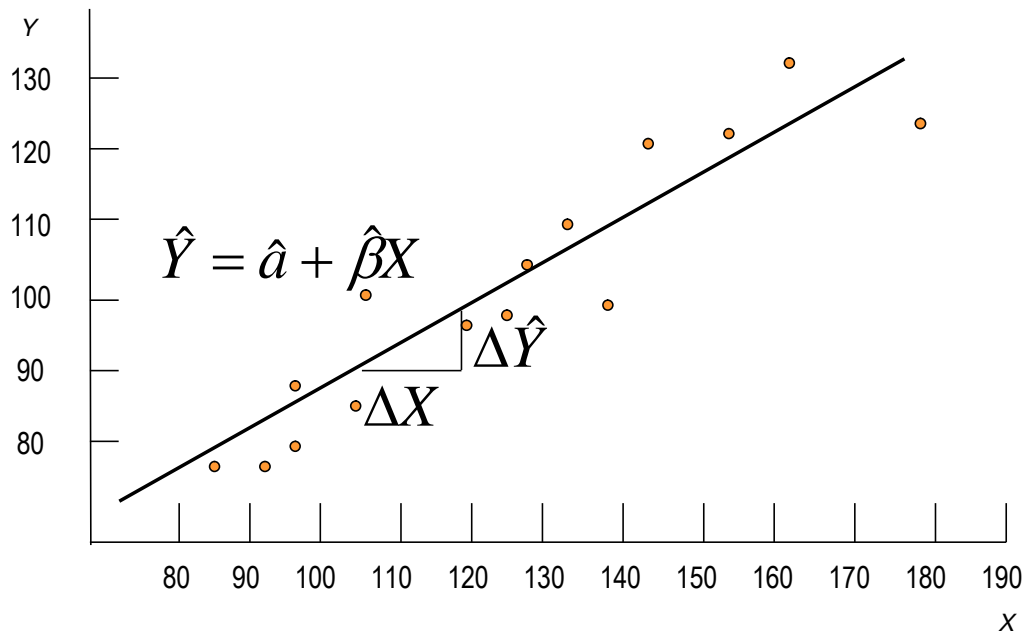
- a
- An intercepted segment of a line
- The point at which a regression line intercepts the Y-axis

Slope

- b
- The inclination of a regression line as compared to a base line
- Rise over run
- Δ - notation for “a change in”

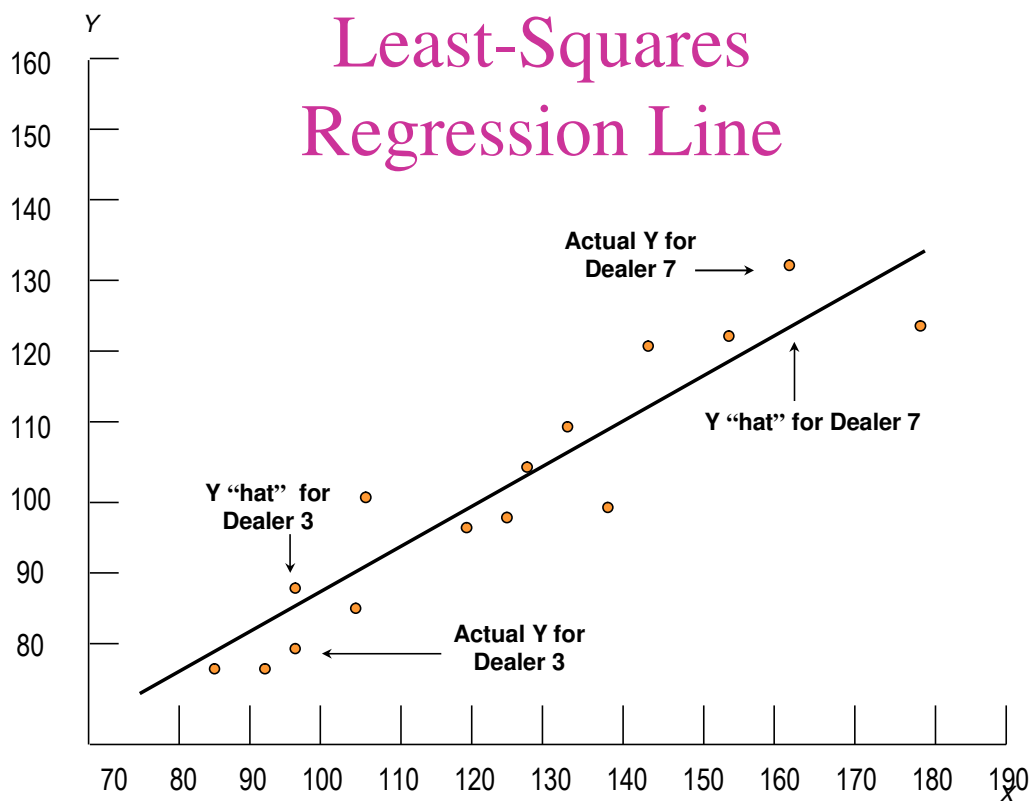


Regression Line and Slope

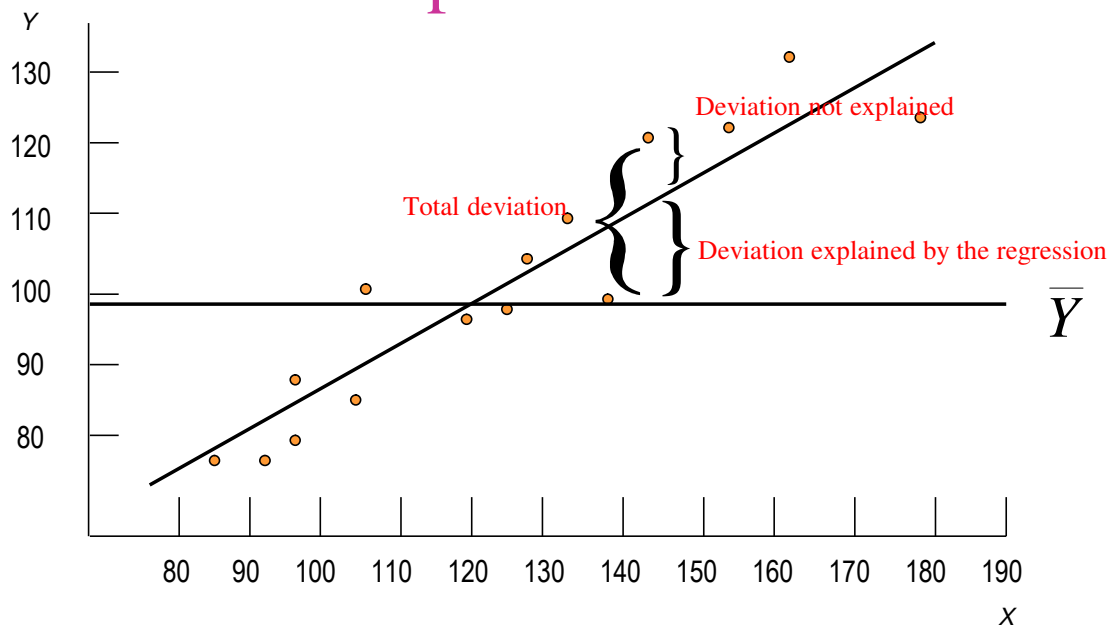


Health Economics Research Method 2003/2.

Least-Squares Regression Line



Scatter Diagram of Explained and Unexplained Variation



Health Economics Research Method 2003/2.

The Least-Square Method

- Uses the criterion of attempting to make the least amount of total error in prediction of Y from X. More technically, the procedure used in the least-squares method generates a straight line that minimizes the sum of squared deviations of the actual values from this predicted regression line.

The Least-Square Method

- A relatively simple mathematical technique that ensures that the straight line will most closely represent the relationship between X and Y.

Regression - Least-Square Method

$$\sum_{i=1}^n e_i^2 \text{ is minimum}$$

$$e_i = Y_i - \hat{Y}_i \quad (\text{The "residual"})$$

Y_i = actual value of the dependent variable

\hat{Y}_i = estimated value of the dependent variable (Y hat)

n = number of observations

i = number of the observation

The Logic behind the Least-Squares Technique

- No straight line can completely represent every dot in the scatter diagram
- There will be a discrepancy between most of the actual scores (each dot) and the predicted score
- Uses the criterion of attempting to make the least amount of total error in prediction of Y from X

Bivariate Regression

$$\hat{a} = \bar{Y} - \hat{\beta}\bar{X}$$

Bivariate Regression

$$\hat{\beta} = \frac{n(\sum XY) - (\sum X)(\sum Y)}{n(\sum X^2) - (\sum X)^2}$$

$\hat{\beta}$ = estimated slope of the line (the “regression coefficient”)

\hat{a} = estimated intercept of the y axis

Y = dependent variable

\bar{Y} = mean of the dependent variable

X = independent variable

\bar{X} = mean of the independent variable

n = number of observations

$$\begin{aligned}\hat{\beta} &= \frac{15(193,345) - 2,806,875}{15(245,759) - 3,515,625} \\ &= \frac{2,900,175 - 2,806,875}{3,686,385 - 3,515,625} \\ &= \frac{93,300}{170,760} = .54638\end{aligned}$$

$$\begin{aligned}\hat{a} &= 99.8 - .54638(125) \\ &= 99.8 - 68.3 \\ &= 31.5\end{aligned}$$

$$\begin{aligned}\hat{a} &= 99.8 - .54638(125) \\ &= 99.8 - 68.3 \\ &= 31.5\end{aligned}$$

$$\begin{aligned}\hat{Y} &= 31.5 + .546(X) \\ &= 31.5 + .546(89) \\ &= 31.5 + 48.6 \\ &= 80.1\end{aligned}$$

$$\begin{aligned}\hat{Y} &= 31.5 + .546(X) \\ &= 31.5 + .546(89) \\ &= 31.5 + 48.6 \\ &= 80.1\end{aligned}$$

Dealer 7 (Actual Y value = 129)

$$\begin{aligned}\hat{Y}_7 &= 31.5 + .546(165) \\ &= 121.6\end{aligned}$$

Dealer 3 (Actual Y value = 80)

$$\begin{aligned}\hat{Y}_3 &= 31.5 + .546(95) \\ &= 83.4\end{aligned}$$

$$\begin{aligned}e_i &= Y_9 - \hat{Y}_9 \\ &= 97 - 96.5 \\ &= 0.5\end{aligned}$$

Dealer 7 (Actual Y value = 129)

$$\begin{aligned}\hat{Y}_7 &= 31.5 + .546(165) \\ &= 121.6\end{aligned}$$

Dealer 3 (Actual Y value = 80)

$$\begin{aligned}\hat{Y}_3 &= 31.5 + .546(95) \\ &= 83.4\end{aligned}$$

$$\begin{aligned}e_i &= Y_9 - \hat{Y}_9 \\ &= 97 - 96.5 \\ &= 0.5\end{aligned}$$

$$\hat{Y}_9 = 31.5 + .546(119)$$

F-Test (Regression)

- A procedure to determine whether there is more variability explained by the regression or unexplained by the regression.
- Analysis of variance summary table

Total Deviation can be Partitioned into Two Parts

- Total deviation equals
- Deviation explained by the regression plus
- Deviation unexplained by the regression

“We are always acting on what has just finished happening. It happened at least 1/30th of a second ago. We think we’re in the present, but we aren’t. The present we know is only a movie of the past.”

Tom Wolfe in
The Electric Kool-Aid Acid Test

Partitioning the Variance

$$(Y_i - \bar{Y}) = (\hat{Y}_i - \bar{Y}) + (Y_i - \hat{Y}_i)$$

Total deviation = Deviation explained by the regression + Deviation unexplained by the regression (Residual error)

\bar{Y} = Mean of the total group

\hat{Y} = Value predicted with regression equation

Y_i = Actual value

$$\sum (Y_i - \bar{Y})^2 = \sum (\hat{Y}_i - \bar{Y})^2 + \sum (Y_i - \hat{Y}_i)^2$$

Total
variation
explained = Explained
variation + Unexplained
variation
(residual)

Sum of Squares

$$SS_t = SS_r + SS_e$$

Coefficient of Determination

$$r^2$$

- The proportion of variance in Y that is explained by X (or vice versa)
- A measure obtained by squaring the correlation coefficient; that proportion of the total variance of a variable that is accounted for by knowing the value of another variable

Coefficient of Determination

$$r^2$$

$$r^2 = \frac{SSr}{SS_t} = 1 - \frac{SSe}{SS_t}$$

Source of Variation

- Explained by Regression
- Degrees of Freedom
 - $k-1$ where k = number of estimated constants (variables)
- Sum of Squares
 - SSr
- Mean Squared
 - $SSr/k-1$

Source of Variation

- Unexplained by Regression
- Degrees of Freedom
 - $n-k$ where n =number of observations
- Sum of Squares
 - SSe
- Mean Squared
 - $SSe/n-k$

r^2 in the Example

$$r^2 = \frac{3,398.49}{3,882.4} = .875$$

Multiple Regression

- Extension of Bivariate Regression
- Multidimensional when three or more variables are involved
- Simultaneously investigates the effect of two or more variables on a single dependent variable
- Discussed in Chapter 24

Microsoft Excel - Correlation Regression Trade Area

File Edit View Insert Format Tools Data Window Help

Arial 10 B I U \$ % , +.0 -.0

C:\Exploring Marketing Research\data\8ed Excel\Correlation Regre

G25 =

	B	C	D	E	F	G
1						
2	MAJOR CITY	POPULATION	RETAIL SALES (000)			
3	Blountstown	13,017	\$108,126			
4	Apalachicola	11,057	\$95,332			
5	Quincy	45,087	\$266,399		POPULATION	RETAIL SALES
6	Monticello	12,902	\$82,837	POPULATION	1	
7	Bristol	7,021	\$10,366	RETAIL SALES	0.846899978	1
8	Madison	18,733	\$103,993			
9	Perry	19,256	\$129,649			
10	Crawfordville	22,863	\$100,849			
11	Quitman	16,450	\$50,529			
12	Bainbridge	28,240	\$302,444			
13	Cairo	23,659	\$166,420			
14	Thomasville	42,737	\$560,412			
15						
16						
17						
18						
19						
20						
21						

TRADE AREA SURROUNDING AREA Regression output Sheet4

Ready CAPS NUM

Microsoft Excel - Correlation Regression Trade Area

File Edit View Insert Format Tools Data Window Help

Arial 10 B I U \$ % , +.0 -.0

C:\Exploring Marketing Research\data\8ed Excel\Correlation Regre

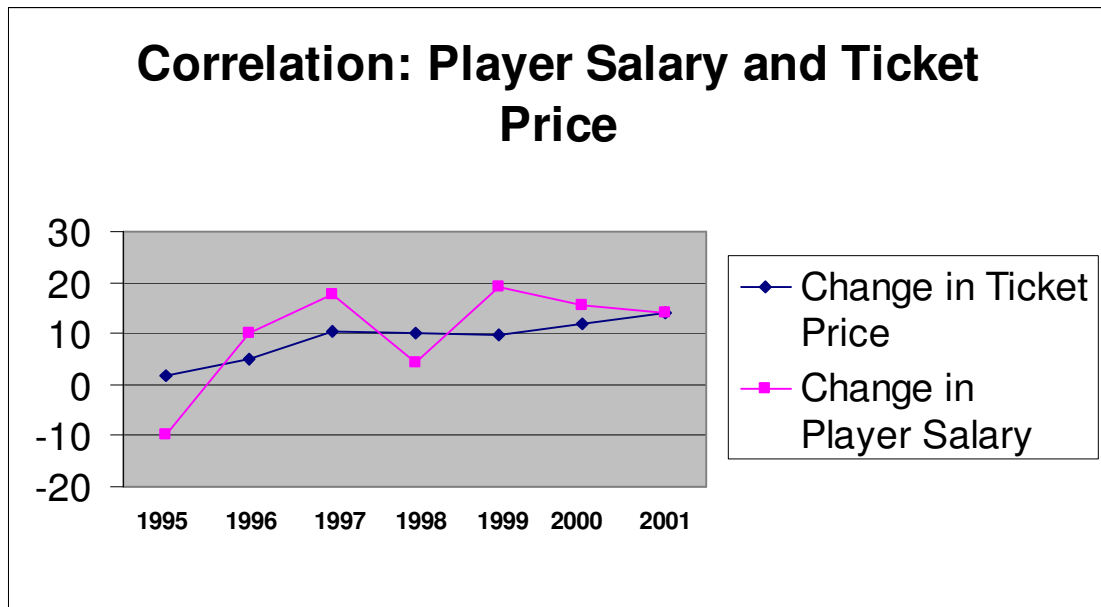
J24 =

	A	B	C	D	E	F	G	H	I	J
1	SUMMARY OUTPUT									
2										
3	<i>Regression Statistics</i>									
4	Multiple R	0.8469								
5	R Square	0.71724								
6	Adjusted R Square	0.688964								
7	Standard Error	83481.02								
8	Observations	12								
9										
10	<i>ANOVA</i>									
11		<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>				
12	Regression	1	1.77E+11	1.77E+11	25.36563	0.00050929				
13	Residual	10	6.97E+10	6.97E+09						
14	Total	11	2.46E+11							
15										
16		<i>Coefficient</i>	<i>Standard Err</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>	<i>Lower 95.0%</i>	<i>Upper 95.0%</i>	
17	Intercept	-66672	51890.89	-1.28485	0.227807	-182292.1167	48948.14378	-182292.1167	48948.14378	
18	POPULATION	10.64056	2.112718	5.03643	0.000509	5.933127658	15.34798917	5.933127658	15.34798917	
19										
20										
21										

TRADE AREA SURROUNDING AREA Regression output Sheet4

Ready CAPS NUM

Correlation Coefficient, $r = .75$



Microsoft Excel - Correlation Regression Trade Area

File Edit View Insert Format Tools Data Window Help

Arial 10 B I U \$ % , +.0 -.0

C:\Exploring Marketing Research\data\8ed Excel\Correlation Regre

Year	Average Ticket	% change	Player Salary	% change	Average Ticket	Player Salary
1994	10.45		1,188,679		1	
1995	10.65	1.9	1,071,029	-9.9	0.991510451	1
1996	11.2	5.1	1,176,967	9.9		
1997	12.36	10.4	1,383,578	17.6		
1998	13.59	10	1,441,406	4.2		
1999	14.91	9.7	1,720,050	19.3		
2000	16.67	11.9	1,988,034	15.6	1	
2001	18.99	13.9	2,264,403	13.9	0.750758207	1

TRADE AREA / SURROUNDING AREA / Regression output / baseball / She

Ready NUM