CHAPTER PROBLEM

Can we predict the cost of subway fare from the price of a slice of pizza?
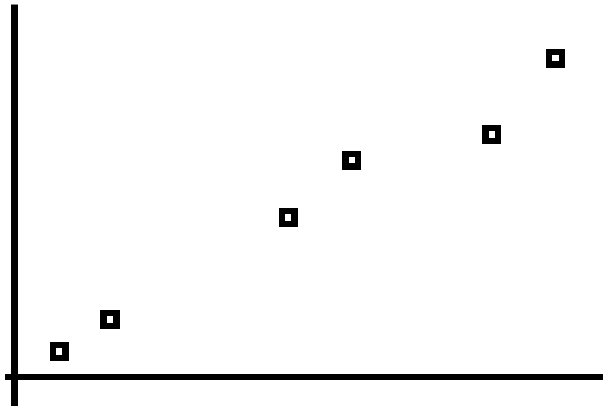
In 1964, Eric Bram, a typical New York City teenager, noticed that the cost of a slice of cheese pizza was the same as the cost of a subway ride. Over the years, he noticed that those two costs seemed to increase by about the same amounts. In 1980, when the cost of a slice of pizza increased, he told the *New York Times* that the cost of subway fare would increase. His prediction proved to be correct.

In the recent *New York Times* article "Will Subway Fares Rise? Check at Your Pizza Place," reporter Clyde Haberman wrote that in New York City, the subway fare and the cost of a slice of pizza "have run remarkably parallel for decades." A random sample of costs (in dollars) of pizza and subway fares are listed in the table below. The table also includes values of the Consumer Price Index (CPI) for the New York metropolitan region, with the index of 100 assigned to the base period from 1982 to 1984. The CPI reflects the cost of a standard collection of goods and services, including such items as a gallon of milk and a loaf of bread. From the table, we see that the paired pizza/subway fare costs are approximately the same for the given years. As a first step, we should examine the data visually. Recall from Section 2-4 that a scatterplot is a plot of (*x,y*) paired data. The pattern of the plotted data points is often helpful in determining whether there is a **correlation**, or association, between the two variables. The scatterplot shown suggests that there is a correlation between the cost of a slice of pizza and the cost of a subway fare. Because an informal conclusion based on an inspection of the scatterplot is largely subjective, we must use other tools for addressing questions such as:

π  If there is a correlation between two variables, how can it be described? Is there an **equation** that can be used to predict the cost of a subway fare given the cost of a slice of pizza?

π  If we can predict the cost of a subway fare, how accurate is that prediction likely to be?

π  Is there also a correlation between the CPI and the cost of a subway fare, and if so, is the CPI better for predicting the cost of a subway fare?

| Year | 1960 | 1973 | 1986 | 1995 | 2002 | 2003 |
|------|------|------|------|------|------|------|
| Cost of Pizza | 0.15 | 0.35 | 1.00 | 1.25 | 1.75 | 2.00 |
| Subway Fare | 0.15 | 0.35 | 1.00 | 1.35 | 1.50 | 2.00 |
| CPI | 30.2 | 48.3 | 112.3 | 162.2 | 191.9 | 197.8 |

| MATH 103 CHAPTER 10 HOMEWORK |
|---|
| **10.2** 1-5, 10, 14, 18, 19, 20, 23, 27 |
| **10.3** 1-5, 9, 10, 18, 19, 20, 23, 27 |

## 10.1  REVIEW AND PREVIEW

In Chapter 9 we presented methods for making _____

from _____ samples. In Section 9-4 we considered two _____

samples, with each value of one sample somehow _____ with a value

from the other sample. In Section 9-4 we considered the _____

between the _____ values, and we illustrated the use of

_____ tests for _____ about the _____

of _____. We also illustrated the _____ of

_____ interval _____ of the _____ of

all such differences. In this chapter, we again consider _____

sample data, but the objective is fundamentally different. In this chapter

we introduce methods for determining whether a _____, or

_____, between two variables exists, and whether

the _____ is _____. For _____

_____ we can identify an _____ that best

_____ the _____ and we can use that equation to

_____ the _____ of one _____

given the value of the other variable.

10.2  CORRELATION
        Key Concept…
        In Part 1 of this section we introduce the _____

        _____ _____ _____, which is a

        _____ measure of the _____ of the

        _____ between _____ variables representing

        _____ data. Using _____ sample data

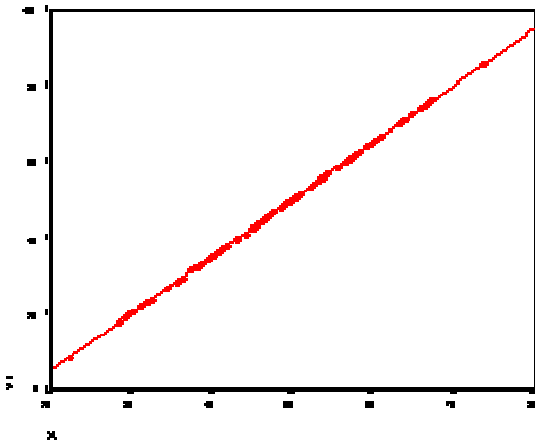        (sometimes called _____ _____), we find the

value of _____, then we use that value to _____ that

there is (or is not) a _____ _____

between the _____ variables. In this section we consider only

_____ relationships, which means that when

_____, the points _____ a

_____-_____ pattern. In Part 2, we discuss

methods of _____ testing for _____.
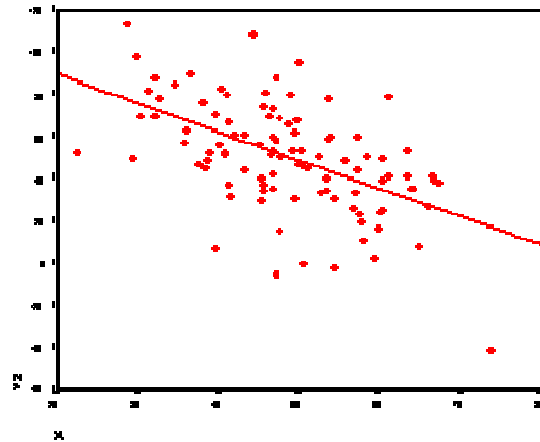
## DEFINITION

A **correlation** exists between two _____ when the _____

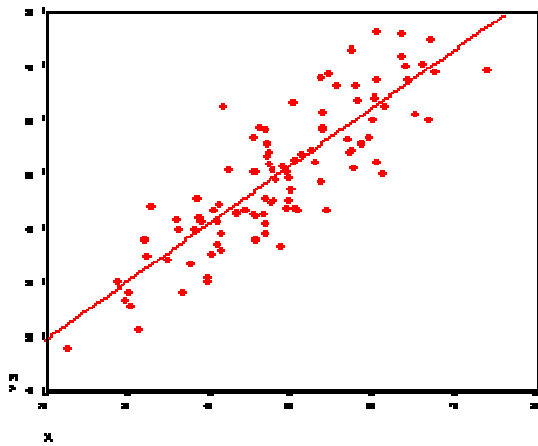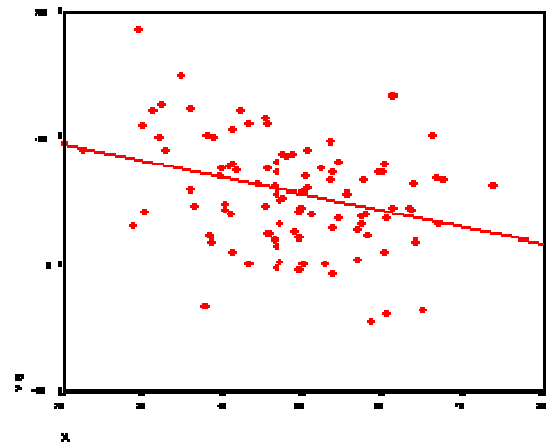of one variable are somehow _____ with the values of the other

variable.

## EXPLORING THE DATA
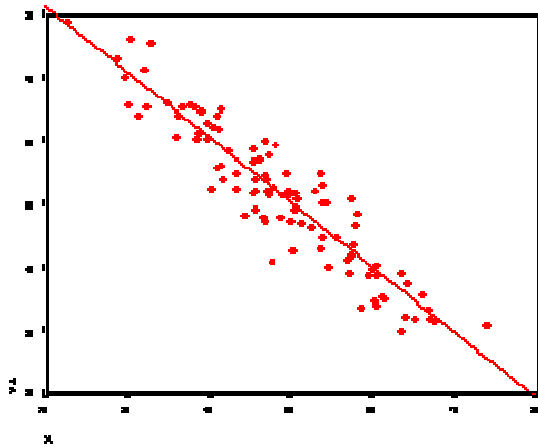
r = 1.00

r = -.54

r = .85



r = -.33



r = -.94



r = .17



r = .42



r = .39

## DEFINITION

The **linear correlation coefficient** *r* measures the _____ of the

_____ _____ between the _____

_____ _____ and _____ in a _____. The

linear correlation coefficient is sometimes referred to as the

_____ _____ _____

_____ _____ in honor of Karl Pearson who

originally developed it.

Because the linear _____ coefficient _____ is calculated using

_____ data, it is a _____ _____. If

we had every pair of _____ values, it would be represented by

_____ (Greek letter rho).

## OBJECTIVE

## NOTATION FOR THE LINEAR CORRELATION COEFFICIENT

$n =$

$\left(\Sigma x\right)^2 =$

$\Sigma =$

$\Sigma xy =$

$\Sigma x =$

$r =$

$\Sigma x^2 =$

$\rho =$

## REQUIREMENTS

1.  The _____ of _____ _____ data is a SRS of

    _____ data.

2.  Visual examination of the _____ must confirm that the

    points _____ a straight-line _____.

3.  Because results can be _____ affected by the presence of

    _____, any _____ must be

    _____ if they are known to be _____. The effects

    of any other _____ should be considered by calculating _____ with

    and without the _____ included.

## FORMULAS FOR CALCULATING $r$

$r =$ —————————————————

$r =$ —————————————————

where _____ is the _____ _____ for the sample value _____ and _____ is

the _____ for the sample value _____.

## INTERPRETING THE LINEAR CORRELATION COEFFICIENT $r$

Computer Software

If the _____ computed from _____ is less than or equal to the

_____ _____, conclude that there is a _____

correlation. Otherwise, there is not _____ evidence to support

the _____ of linear _____.

Table A-5

If the _____ _____ of _____, denoted _____, exceeds

the value in Table A-5, conclude that there is a _____ correlation.

Otherwise, there is not sufficient evidence to _____ the conclusion
of a linear correlation.

## ROUNDING THE LINEAR CORRELATION COEFFICIENT $r$

Round the _____ _____ _____ _____ to

_____ decimal places so that its value can be compared to critical values in

Table A-5. Keep as many decimal places during the process and then _____

at the end.

## PROPERTIES OF THE LINEAR CORRELATION COEFFICIENT $r$

1.  The value of _____ is always between _____ and _____ inclusive. That is

    _____.

2.  If all values of _____ variable are _____ to a

    different _____, the value of _____ _____ _____

    change.

3.  The value of _____ is _____ affected by the choice of _____ or _____.

4.  _____ measures the _____ of a _____ relationship.

    It is not designed to measure the strength of a _____ that

is _____ linear.

5. _____ is very sensitive to _____ in the sense that a _____

outlier can _____ affect its value.

## COMMON ERRORS INVOLVING CORRELATION

1. A common _____ is to _____ that

_____ implies _____.

2. Another error arises with data based on _____. Average

_____ _____ variation and may _____

the _____ _____.

3. A third error involves the property of _____. If there is no

linear _____, there might be some other _____

that is not _____.

## PART 2: FORMAL HYPOTHESIS TEST

| HYPOTHESIS TEST FOR CORRELATION (USING TEST STATISTIC $r$) |
|---|
| **NOTATION** |
| $n =$                                            $\rho =$<br><br><br>$r =$ |

**HYPOTHESES**

**TEST STATISTIC:** *r*

Critical values: Refer to Table _____

**CONCLUSION**

Example 1: The heights and weights of a sample of 9 supermodels were measured. Using a TI-83/84 Plus calculator, the linear correlation coefficient of the 9 pairs of measurements is found to be 0.360. Is there sufficient evidence at the 5% level to support the claim that there is a linear correlation between the heights and weights of supermodels? Explain.

## HYPOTHESIS TEST FOR CORRELATION (USING P-VALUE FROM A t-TEST)

### HYPOTHESES

### TEST STATISTIC

$$t = \frac{\rule{4cm}{0.4pt}}{}$$

P-value: Use _____ or Table _____ with _____ degrees of freedom.

### CONCLUSION

Example 2: The paired values of the CPI and the cost of a slice of pizza are listed below.

| CPI | 30.2 | 48.3 | 112.3 | 162.2 | 191.9 | 197.8 |
|---|---|---|---|---|---|---|
| Cost of Pizza | 0.15 | 0.35 | 1.00 | 1.25 | 1.75 | 2.00 |

a. Construct a scatterplot

b. Find the value of the linear correlation coefficient $r$

c. Find the critical values of $r$ from Table A-5 using a significance level of 0.05.

d. Determine whether there is sufficient evidence to support a claim of a linear correlation between the two variables.

10.3  INFERENCES ABOUT TWO MEANS: INDEPENDENT SAMPLES
Key Concept…
In section 10-2, we presented methods for finding the value of the

_____ correlation _____ _____ and for

determining whether there is a _____ correlation between two variables. In Part 1 of this section, we find the _____ of the _____ line that _____ fits the _____ sample data. The equation algebraically describes the _____ between the two variables. The best-fitting straight line is called the _____ line, and its equation is called the _____ equation. We also present methods for using the regression equation to make _____ . In Part 2 we discuss _____ change, _____ points, and _____ plots as a tool for _____ correlation and _____ results.

## PART 1: BASIC CONCEPTS OF REGRESSION

Two variables are sometimes related in a _____ way, meaning that given a value for one variable, the _____ of the other variable is _____ determined without any _____, as in the equation $y = 6x + 5$. Statistics courses focus on _____ models, which are equations with a variable that is _____ completely by the other variable.

**DEFINITION**

Given a collection of _____ sample data, the **<u>regression equation</u>**

algebraically describes the _____ between the two variables

_____ and _____. The _____ of the _____ equation is

called the **<u>regression line (aka line of best fit, or least-squares line)</u>**.

The regression equation expresses a relationship between the

_____ variable _____ (aka _____ variable or

_____ variable) and _____ (called the _____

variable, or _____ variable). The slope and $y$-intercept can be

found using the following formulas:

**The _____ line fits the _____ points _____!**

## OBJECTIVE

## NOTATION

|  | Population Parameter | Sample Statistic |
|---|---|---|
| y-intercept of regression equation |  |  |
| Slope of regression equation |  |  |
| Equation of the regression line |  |  |

## REQUIREMENTS

1. The _____ of _____ _____data is a SRS of

_____ data.

2. Visual examination of the _____ must confirm that the

points _____ a straight-line _____.

3. Because results can be _____ affected by the presence of

_____, any _____ must be

_____ if they are known to be _____. The effects

of any other _____ should be considered by calculating _____ with

and without the _____ included.

**FORMULAS FOR FINDING THE SLOPE \_\_\_\_ AND *y*-INTERCEPT \_\_\_\_ IN**

**THE REGRESSION EQUATION _____**

Slope:

where \_\_\_\_ is the _____ correlation coefficient, \_\_\_\_\_ is the

_____ _____ of the \_\_\_\_ values, and \_\_\_\_\_ is the

standard deviation of the \_\_\_\_\_ values

*y*-intercept:

**ROUNDING THE SLOPE AND THE *y*-INTERCEPT**

Round \_\_\_\_\_ and \_\_\_\_\_ to _____ _____ _____.

**USING THE REGRESSION EQUATION FOR PREDICTIONS**
1. Use the regression equation for _____ only if the

    _____ of the _____ line on the _____

    confirms that the _____ _____ _____ the points
    reasonably well.

2. Use the regression equation for _____ only if the _____

correlation coefficient _____ indicates that there is a _____

correlation between the two variables.

3. Use the regression line for predictions only if the _____ do not go

much _____ the _____ of the available _____
data.

4. If the regression equation does not appear to be _____ for

making _____, the best _____ value is its

_____ estimate, which is its _____ _____.

Example 1: The paired values of the CPI and the cost of a slice of pizza are listed below.

| CPI | 30.2 | 48.3 | 112.3 | 162.2 | 191.9 | 197.8 |
|---|---|---|---|---|---|---|
| Cost of Pizza | 0.15 | 0.35 | 1.00 | 1.25 | 1.75 | 2.00 |

a. Find the regression equation, letting the first variable be the predictor ($x$) variable.

b. Find the best predicted cost of a slice of pizza when the CPI is 182.5.

## PART 2: BEYOND THE BASICS OF REGRESSION

### DEFINITION

In working with two variables _____ by a regression equation, the

**marginal change** in a _____ is the _____ that it

changes when the other variable changes by exactly _____ unit. The slope _____

in the regression equation represents the _____ _____ in _____

that occurs when _____ changes by _____ unit.

### DEFINITION

In a _____, an **outlier** is a point lying _____ away from the

other data points. Paired sample data may include one or more **influential points**,

which are _____ that _____ affect the _____ of the

_____ _____.

## DEFINITION

For a pair of sample _____ and _____ values, the **residual** is the _____

between the _____ sample value of _____ and the _____ that is

_____ by using the _____ equation. That is,

Residual = _____ ___ - _____ ___ = _____

## DEFINITION

A _____ line satisfies the **least-squares property** if the _____ of

their _____ is the _____ sum possible.

## DEFINITION

A **residual plot** is a _____ of the _____ values after each

of the _____ values has been _____ by the

_____ value _____. That is, a residual plot is a graph of

the points _____.