Disadvantages over Survey Method:

things such as awareness, beliefs, feelings and preferences cannot be observed

the observed behavior patterns can be rare or too unpredictable thus increasing the data collection costs and time requirements

- 3. Experimental method - a method designed for collecting data under controlled conditions. An experiment is an operation where there is actual human interference with the conditions that can affect the variable under study. This is an excellent method of collecting data for *causation studies*. If properly designed and executed, experiments will reveal with a good deal of accuracy, the effect of a change in one variable on another variable.
- 4. Use of existing studies e.g., census, health statistics, and weather bureau reports

Two types:



archers who field sourd ave one studies on the area of interest rectly for information needed

5. Registration method e.g., car registration, student registration, and hospital admission

General Classification of Collecting Data

Definition. **Census** or **complete enumeration** is the process of gathering information from every unit in the population.

- not always possible to get timely, accurate and economical data
- costly, especially if the number of units in the population is too large

Definition. Survey sampling is the process of obtaining information from the units in the selected sample.

Advantages of Survey Sampling:

- reduced cost
- greater speed

PRESENTATION OF DATA Chapter 3

Textual Presentation

data incorporated to a paragraph of text

Example

At last count, 38 airlines were operating Boeing 707's, 720's, and 727's over the world's airlines. The far-flung Boeing fleet has now logged an estimated 1,803,704,000 miles (22,855,948,000 kms.) and has massed approximately 4,096,000 revenue flight hours. Passenger totals stand at upwards of 71.6 million.

Advantages

- This presentation gives emphasis to significant figures and comparisons
- It is simplest and most appropriate approach when there are only a few numbers to be

Disadvantages

- r paragraph, the presentation
- When a large mass of quantitative of the and the method in a flat or name becomes almost incomprehensible Paragraphs can be provided to read me to read especially the same words are repeated so many times

Tabular Presentation

• the systematic organization of data in rows and columns

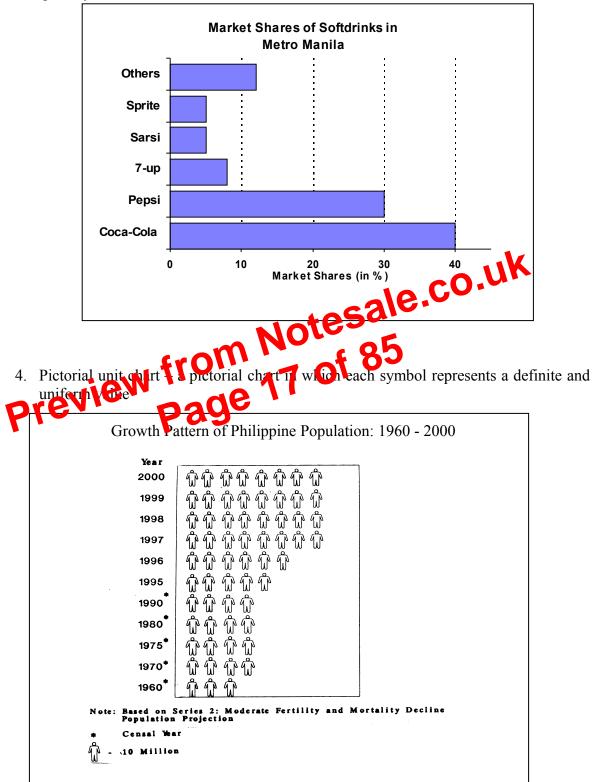
Advantages

- more concise than textual presentation •
- easier to understand
- facilitates comparisons and analysis of relationship among different categories
- presents data in greater detail than a graph

Parts of a Formal Statistical Table

1. Heading consists of a *table number*, *title*, and *headnote*. The title is a brief statement of the nature, classification and time reference of the information presented and the area to which the statistics refer. The headnote is a statement enclosed in

3. Bar Chart - consists of a series of rectangular bars where the length of the bar represents the quantity or frequency for each category if the bars are arranged horizontally. If the bars are arranged vertically, the height of the bar represents the quantity.



Pearson's First and Second Coefficients of Skewness

Example: Refer to the final grade of 110 Soc Sci 101 students.

$$\overline{X}$$
 = 74.1 Md = 75 Mo = 84 s = 11.25

Using the first formula,

$$Sk = \frac{74.1 - 84}{11.25} = -0.88$$

Using the second formula,

$$Sk = \frac{3(74.1 - 75)}{11.25} = -0.24$$



THE PROBABILITY CONCEPT AND SOME PROPERTIES

Probability analysis is based on the following simple postulates.

- Postulate 1. $0 \le P(O_i) \le 1$ for any simple event O_i
- Postulate 2. The probability for any event E is the sum of the probability of the simple events that constitute E.
- Postulate 3. P(S) = 1, where S is the sample space, and $P(\phi) = 0$, where ϕ is the null space.

Approaches to Assigning Probabilities

1. A Priori or Classical Probability – probability is determined even before the experiment is performed using the following rule:

If an experiment can result in any one of N different equally likely outcomes, and if exactly n of these outcomes correspond to event A, then the probability of event A is

$$P(A) = \frac{\text{no. of sample points in A}}{\text{no. of sample points in S}} = \frac{n}{N}$$
2. A Posteriori or Relative Trequency or Empirical Probability - probability is determined by repeating the experiment a large number of times using the following rule:

$$P(A) = \frac{P(A)}{P(A)} = \frac{P(A) - P(A)}{P(A)}$$

3. **Subjective Probability** – probability is determined by the use of intuition, personal beliefs, and other indirect information.

Examples:

- 1. Find the errors in each of the following statements:
 - a. The probability that it will rain tomorrow is 0.40 and the probability that it will not rain tomorrow is 0.52.
 - b. The probabilities that a printer will make 0, 1, 2, 3, or 4 or more mistakes in printing a document are, respectively, 0.19, 0.34, -0.25, 0.43, and 0.29.
 - c. The probabilities that an automobile salesperson will sell 0, 1, 2, or 3 cars on any given day in February are, respectively, 0.19, 0.38, 0.29, and 0.15.
 - d. On a single draw from a deck of playing cards the probability of selecting a heart is 1/4, the probability of selecting a black card is 1/2, and the probability of selecting both a heart and a black card is 1/8.
- 2. Answer the following:
 - a. In tossing a fair coin, what is the probability of getting a head? Of either a head or tail? Of neither a head nor tail?
 - b. In tossing a fair die, what is the probability of getting a 3? Of getting an even number? Of getting a number greater than 6?

If an operation can be performation n Theorem. vays, and for each of these a second by can be performed an_2 ways, then the two operations can be performed in m₁n₂ wavs.

Example:

- 1. How many sample points are there in the sample space when a pair of balanced dice is thrown once?
- 2. If a travel agency offers special weekend trips to 12 different cities, by air, rail, or bus, in how many ways can such a trip be arranged?
- Theorem. (Multiplication Rule) If an operation can be performed in n_1 ways, if for each of these a second operation can be performed in n_2 ways, if for each of the first two a third operation can be performed in n_3 ways, and so on, then the sequence of k operations can be performed in $n_1n_2 \dots n_k$ ways.

Examples:

- 1. How many even three-digit numbers can be formed from the digits 3, 5, 6, 7, 9, 2 if
 - a. all digits are distinct?
 - b. numbers must be divisible by 5 and repetition is not allowed?
 - c. Numbers must be even and repetition is not allowed?

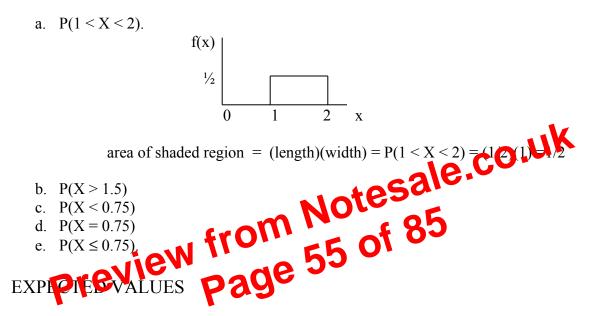
2. The probability density function can not be represented in tabular form.

Example:

A continuous random variable X that can assume values between 0 and 2 has a density function given by

$$f(x) = \begin{cases} 0.5 & \text{for } 0 < x < 2\\ 0 & \text{otherwise} \end{cases}$$

Find the following probabilities:



Definition. Let X be a discrete random variable with probability distribution

The mean or expected value of X is

$$\mu = E(X) = \sum_{i=1}^{n} x_i f(x_i)$$

Examples:

1. Find the mean of the random variables X and Y of Experiment No. 1.

E(Y) = (-3)(1/8) + (-1)(3/8) + (1)(3/8) + (3)(1/8) = 0

- 2. Find the expected number of correct matches in Experiment No. 2.
- 3. In a gambling game a man is paid P50 if he gets all heads or all tails when 3 coins are tossed, and he pays out P30 if either 1 or 2 heads show. What is his expected gain?

Theorem. Let X be a discrete random variable with probability distribution

The mean or expected value of the random variable g(X) is

$$E(g(X)) = \sum_{i=1}^{n} g(x_i) f(x_i)$$

- Example: A used car dealer finds that in any day, the probability of selling no car is 0.4, one car is 0.2, two cars is 0.15, 3 cars is 0.10, 4 cars is 0.08, five cars is 0.06 and six cars is 0.01. Let g(X) = 500 + 1500X represent the salesman's daily earnings, where X is the number of cars sold. Find the salesman's expected duty earnings.
- Definition. Let X be a random variable with meaning the Othe variance of X is

$$\sigma^{2} = Var(X) (A - \mu)^{2}$$
Definition. Let V a discrete rand μ variable with probability distribution
$$\frac{X}{P(X=x)} | \begin{array}{c} x_{1} & x_{2} & \dots & x_{n} \\ f(x_{1}) & f(x_{2}) & \dots & f(x_{n}) \end{array}$$

The variance of X is

$$\sigma^{2} = Var(X) = E(X - \mu)^{2} = \sum_{i=1}^{n} (x_{i} - \mu)^{2} f(x_{i})$$

Theorem. Computational Formula for σ^2

$$Var(X) = E(X^2) - [E(X)]^2$$

Example : In Experiment No. 1, find the variance of X..

Using the definition of Var(X),
E(X) = 1.5
Var(X) =
$$\sum_{i=1}^{4} (x_i - 1.5)^2 f(x_i)$$

= $(0-1.5)^2(1/8) + (1-1.5)^2(3/8) + (2-1.5)^2(3/8) + (3-1.5)^2(1/8) = 0.75$

Using the computational formula of the Var(X),

 $Var(X) = E(X^2) - [E(X)]^2 = 3 - (1.5)^2 = 0.75$ Properties of the Mean and Variance

Let X and Y be random variables (discrete or continuous) and let **a** and **b** be constants.

1. E(aX + b) = a E(X) + b

a

a. if b = 0, then E(aX) = a E(X). b. if a = 0, then E(b) = b.

- 2. E(X+Y) = E(X) + E(Y)E(X-Y) = E(X) - E(Y)
- 3. E(XY) = E(X)E(Y) if X and Y are independent.

Special Cases:

4.
$$E[X - E(X)] = 0.$$

5.
$$\operatorname{Var}(aX + b) = a^2 \operatorname{Var}(X).$$

а

a Cases: a. if b = 0, then $Var(aX) = a^2 Va(X)$. b. if a = 0, then Var(aX) = 0. and independent then A = 57 of 85 Var(X) = Var(X) = 0. Special Cases: 6. If X and Var(X + Y) = Var(X) + Var(Y)Var(X - Y) = Var(X) + Var(Y)

Example :

If X and Y are independent random variables with E(X) = 3, E(Y) = 2, Var(X) = 2 and Var(Y)=1, find

- a. E(3X + 5)
- b. Var(3X+5)
- c. E(XY)
- c. Var(3X 2Y)

- Notation: t_{α} is the t-value leaving an area of α in the right-tail of the t-distribution. That is, if $T \sim t_{(v)}$ then t_{α} is such that $P(T > t_{\alpha}) = \alpha$.
- Since the t-distribution is symmetric about zero, $t_{1-\alpha} = -t_{\alpha}$

Examples:

- 1. Find the following values on the t -table:
 - (a) $t_{0.025}$ when v = 14.
 - (b) $t_{0.99}$ when v=10.
- 2. Find k such that P(k < T < 2.807) = 0.945 when $T \sim t_{(23)}$

3. A manufacturing firm claims that the batteries used in their electronic games will last an average of 30 hours. To maintain this average, 16 batteries are tested each month. If the computed t-value falls between $-t_{0.025}$ and $t_{0.025}$, the firm is satisfied with its claim. What conclusion should the firm draw from a sample that has mean $\overline{X} = 27.5$ hours and standard deviation S = 5 hours? Assume the distribution of battery lives to be approximately normal.

Chapter 12 Estimation of Parameters

Definition. Statistical inference refers to methods by which one uses sample information to make inferences or generalizations about a population.

Two Areas of Statistical Inference

- 1. Estimation
 - point estimation
 - interval estimation
- 2. Hypothesis Testing

BASIC CONCEPTS IN ESTIMATION

Point Estimation

Definition. An estimator is any statistic whose value is used to stille an unknown parameter. A realized value of an estimator in a stimate.

For example, the sample mean \overline{k}_{μ} is an examptor of the population mean μ . marks: 68 0

Remarks:

An estimator is said to be aboved if the average of the estimates it produces under 1. repeated sampling is equal to the true value of the parameter being estimated.

Examples: Under random sampling, the sample mean is an unbiased estimator of the population mean, that is, $E(\overline{X}) = \mu$.

> Under random sampling with replacement, S^2 is an unbiased estimator of σ^2 , but S on the other hand is a biased estimator of σ with the bias becoming insignificant for large samples.

2. A parameter can have more than one unbiased estimator. We would naturally choose the unbiased estimator with the smallest variance.

Interval Estimation

Definition. An **interval estimator** of a population parameter is a rule that tells us how to calculate two numbers based on sample data, forming an interval within which the parameter is expected to lie. This pair of numbers, (a,b), is called an interval estimate or confidence interval.

Examples:

- 1. An electrical firm manufactures light bulbs that have a length of life that is normally distributed, with a standard deviation of 40 hours. If a random sample of 25 bulbs has a mean life of 780 hours, find a 95% confidence interval for the population mean of all bulbs produced by this firm.
- 2. Regular consumption of presweetened cereals contribute to tooth decay, heart disease, and other degenerative diseases, according to a study by Dr. M. Albreight of the National Institute of Health and Dr. D. Solomon, Professor of Nutrition and Dietetics at the University of London. In a random sample of 20 similar servings of Alpha-Bits, the mean sugar content was 11.3 grams with a standard deviation of 2.45 grams. Assuming that the sugar content is normally distributed, construct a 95% confidence interval for the mean sugar content for single servings of Alpha-Bits.
- A random sample of 100 automobile owners shows that an automobile is driven on the average 23,500 kilometers per year, in the state of Virginia, with a standard deviation of 3900 kilometers. Construct a 99% confidence interval for the average number of kilometers an automobile is driven annually in Virginia.

ESTIMATING THE DIFFERENCE BETWEIOLTWO POPULATION MEANS Difference we have two populations versions μ_1 and μ_2 and standard deviations σ_1 and σ_2 ,

respectively, a point estimator of the difference between μ_1 and μ_2 is the statistic $\overline{X}_1 - \overline{X}_2$.

Types of Sampling:

- selecting two independent samples
- paired sampling

Paired sampling is used to overcome the difficulty imposed by extraneous differences between two groups when testing the difference between 2 means. This is achieved by "matching" or studying 2 related samples. Matching may be achieved by:

- using the same subject in the 2 samples
- pairing of subjects with respect to any extraneous variable which might affect or influence the outcome.

Example:

In a random sample of 200 students who enrolled in Math 17, 138 passed on their first take. Construct a 95% confidence interval for the population proportion of students who passed Math 17 on their first take.

ESTIMATING THE DIFFERENCE OF TWO PROPORTIONS

Given 2 independent random samples of size n_1 and n_2 , a point estimator of the difference between the two proportions p_1 and p_2 is given by $\hat{p}_1 - \hat{p}_2 = \frac{X}{n_1} - \frac{Y}{n_2}$, where X is the number of successes in n_1 trials (first sample) and Y is the number of successes in n_2 trials (second sample).

An approximate $(1-\alpha)100\%$ confidence interval for $p_1 - p_2$ when n_1 and n_2 are large is

$$\begin{pmatrix} (\hat{p}_{1} - \hat{p}_{2}) - z_{\alpha/2} \sqrt{\frac{\hat{p}_{1}\hat{q}_{1}}{n_{1}} + \frac{\hat{p}_{2}\hat{q}_{2}}{n_{2}}}, & (\hat{p}_{1} - \hat{p}_{2}) = 0 \\ \hline n_{1} + \frac{\hat{p}_{2}\hat{q}_{2}}{n_{2}} \end{pmatrix}$$

Example:

SAMPLE SIZE DETERMINATION

Sample Size for Estimating μ

In random sampling, if \overline{X} will be used to estimate μ , we can be $(1-\alpha)100\%$ confident that the error will not exceed a specified amount, e, when the sample size is

$$n = \left(\frac{z_{\alpha/2}\sigma}{e}\right)^2$$

Based on 2 independent samples			
Но	Test Statistic	На	Critical region
a. σ_1^2 and σ_2^2 known			
$\mu_1 - \mu_2 = d_o$	$Z = \frac{(\overline{X}_{1} - \overline{X}_{2}) - d_{o}}{\sqrt{(\sigma_{1}^{2}/n_{1}) + (\sigma_{2}^{2}/n_{2})}}$	μ_1 - μ_2 < d _o μ_1 - μ_2 > d _o	$z < - z_{\alpha}$
		$\mu_1 - \mu_2 \neq d_0$	$ z > z_{\alpha/2}$
b. $\sigma_1^2 = \sigma_2^2$ but unknown			
	$t = \frac{(\overline{X}_1 - \overline{X}_2) - d_o}{S_p \sqrt{(1/n_1) + (1/n_2)}}$		
$\mu_1 - \mu_2 = d_o$	$S_p \sqrt{(1/n_1) + (1/n_2)}$	μ_1 - $\mu_2 < d_o$ μ_1 - $\mu_2 > d_o$	$t \leq -t_{lpha} \ t > t_{lpha}$
	$\upsilon = n_1 + n_2 - 2$	$\mu_1 - \mu_2 \neq d_o$ $\mu_1 - \mu_2 \neq d_o$	$ t > t_{\alpha}$
	$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$		
	$n_1 + n_2 - 2$		
c. $\sigma_1^2 \neq \sigma_2^2$ and unknown			
$\mu_1 - \mu_2 = d_0$	$t = \frac{(\bar{X}_1 - \bar{X}_2) - d_o}{\sqrt{(S_1^2/n_1) + (S_2^2/n_2)}}$	μ_1 - μ_2 < d $_{\rm o}$	t < - t
$\mu_1 - \mu_2 \mathbf{u}_0$		$\mu_1 - \mu_2 < d_o$ $\mu_1 - \mu_2 > d_o$	$\mathbf{c} \mathbf{c} \mathbf{c} \mathbf{c} \mathbf{c} \mathbf{c} \mathbf{c} \mathbf{c} $
	$\upsilon = \frac{\left(S_1^2/n_1 + S_2^2/n_2\right)^2}{\left(S_1^2/n_1\right)^2 + \left(S_2^2/n_2\right)^2}$	µ₁ - u - 6	$ t > t_{\alpha/2}$
		OTESE	
	$n_1 - 1$ $n_2 - 1$ on 2 related simples		
• Based (Test Statistic		
Но	Test Statistic	На	Critical region
PIC.	Play	$\mu_{\rm D} < d_{\rm o}$	$t < -t_{lpha}$
$\mu_D = d_o$	S_d/\sqrt{n}	$\mu_{\rm D} > d_{\rm o}$	$t > t_{\alpha}$
	$\upsilon = n - 1$	$\mu_{\rm D} \neq d_{\rm o}$	$ t > t_{\alpha/2}$

Remark: The remarks made in Chapter 8.3 relative to the use of a given statistic apply to the tests described here.

Examples:

- 1. A statistics test was given to 50 girls and 75 boys. The girls made an average of 80 with a standard deviation of 4 and the boys had an average of 86 with a standard deviation of 6. Is there sufficient evidence at 0.05 level of significance that the average grades of girls and boys differ?
- 2. A study was made to determine if the subject matter in a physics course is better understood when a lab constitutes part of the course. Students were allowed to choose between a 3-unit course without lab and a 4-unit course with lab. In the section with lab, a sample of 11 students had an average grade of 85 with a standard deviation of 4.7, and in the section without lab, a sample of 17 students had an average grade of 79 with a standard deviation of