

Introduction to Descriptive and Inferential Statistics

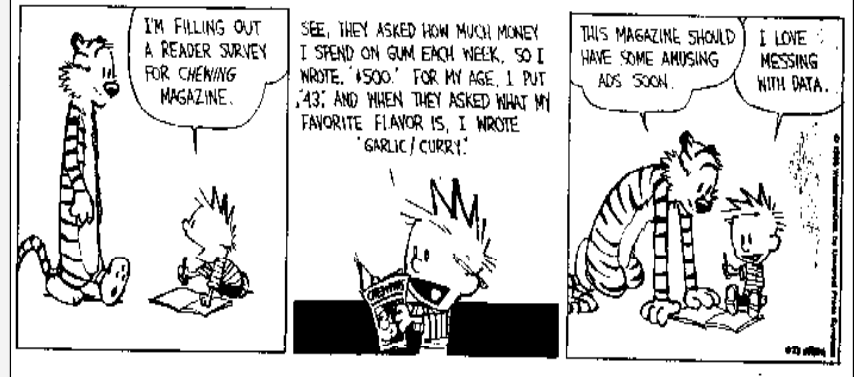


Jaranit Kaewkungwal, Ph.D.
Faculty of Tropical Medicine
Mahidol University

1

Data & Variables

CALVIN and HOBBS



2

Types of Data

QUALITATIVE

Data expressed by type

Data that has been described

QUANTITATIVE

Data classified by numeric value

Data that has been measured or counted



QUALITATIVE and QUANTITATIVE data are not mutually exclusive

Adapted from: Dr. Craig Jackson, University of Central England



Types of Data: Qualitative (Categorical) Data

NOMINAL DATA

- values that the data may have do not have specific order
- values act as labels with no real meaning
- Binomial: two possible values (categories, states)
- Multinomial: more than two possible values (categories, states)

e.g. *Health status* *healthy* = 1 *sick* = 2

e.g. *Treatment* *new regimen* = 1 *standard regimen* = 2

e.g. *hair colour* *brown* = 1 *blond* = 2 *black* = 100

ORDINAL DATA

- values with some kind of ordering
- data that has been measured or counted

e.g. *social class*: *upper* = 1 *middle* = 2 *working* = 3

e.g. *glioblastoma tumor grade*: 1 2 3 4 5

e.g. *position in a race*: 1st 2nd 3rd

Adapted from: Dr. Craig Jackson, University of Central England



Types of Data: Quantitative Data

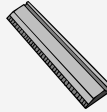
DISCRETE

- distinct or separate parts, with no finite detail
e.g children in family



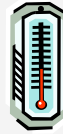
CONTINUOUS

- between any two values, there would be a third
e.g between meters there are centimetres



INTERVAL

- equal intervals between values and an arbitrary zero on the scale
e.g temperature gradient



RATIO

- equal intervals between values and an absolute zero
e.g body mass index



Adapted from: Dr. Craig Jackson, University of Central England

UCE
Birmingham

Examples of Data Coding

Site Number: [][][] Subject Number: [][][][] Chk: [] Initials: [][][][]

Demographics

1. Gender
 1 Male at birth 2 Female at birth
 3 Surgical / medical correction to male
 4 Surgical / medical correction to female

2. Date of Birth
 [][] dd [][] mm [][] yy

3. Race / Ethnicity
 1 Thai 2 Hill-Tribe 3 Unknown 4 Other, specify: _____

Nominal/Cat. Var

15. ในช่วง 6 เดือนที่ผ่านมา ท่านมีเพศสัมพันธ์กับคู่นอนที่ท่านอยู่ด้วยบ่อยแค่ไหน (In the last six months, how often have you had sexual intercourse with your live-in partner? Would you say:)

1 น้อยกว่าเดือนละครั้ง (Less than once a month)
 2 ทุกเดือนแต่ไม่ทุกสัปดาห์ (Every month, but not every week)
 3 ทุกสัปดาห์แต่ไม่ทุกวัน (Every week, but not every day)
 4 เกือบทุกวัน (About every day)
 88 ไม่ทราบ (Don't Know) **Ordinal/Cat. Var**
 99 ไม่ตอบ (Non-response)

Examples of Data Coding

Site Number: [][][] Subject Number: [][][][] Chk: [] Date of Exam: [][] dd [][] mm [][] yy

Physical Exam

1. Height: [][] cm
 2. Weight: [][][] kg
 3. Temperature: [][] °C
 4. Pulse: [][] /min
 5. Respiration: [][] /min
 6. Blood Pressure: [][] / [][] mmHg
 systolic diastolic

Instructions: Physical exam directed by history, systems review and symptoms only.

1 2 99

Body System Normal Abnormal Not Done Comment on all abnormalities

7. HEENT _____
 8. Lymph Nodes _____
 9. Cardiovascular _____

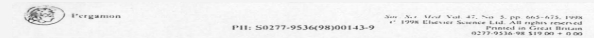
see pp. 4 & 7 Lab 95.6 Page 50

Study #: V0833 Subject Identification
 Drug: MN Recombinant Glycoprotein 120 HIV-1 Vaccine Subject Number: [][][][]
 Study: Phase III Recovering IVU in Thailand Subject Initials: [][][]
 Specimen Date: [][] [][] [][]

CHEMISTRY	Results	ND	HEPATIC			
			NA	Pos	Neg	ND
Creatinine		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
ALT/SGPT		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
WBC		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
RBC		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Hemoglobin		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Hematocrit		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Platelet Count		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Segmented Neutrophils		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Banded Neutrophils		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Lymphocytes		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Monocytes		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Eosinophils		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Basophils		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Other, Specify:		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Urea Nitrogen		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Urine Albumin		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Urine Nitrite		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Urine Esterase		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Urine Glucose		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Urine Pregnancy		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
Urine Dipates		<input type="checkbox"/>		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

If "Not Done", please check ND boxes to indicate this.
 N/A means "Not Applicable" (for example, Pregnancy - males).

Example of Descriptive Statistics



FACTORS RELATED TO PHYSICAL ACTIVITY: A STUDY OF ADOLESCENTS
 BUNAR VILLIALLAESSON* and THOROLFUR THORLINDSSON*
 *Department of Nursing, University of Iceland, Reykjavik 26, 101, Reykjavik, Iceland and *Faculty of Social Sciences, University of Iceland, Sturlugata, 101 Reykjavik, Iceland

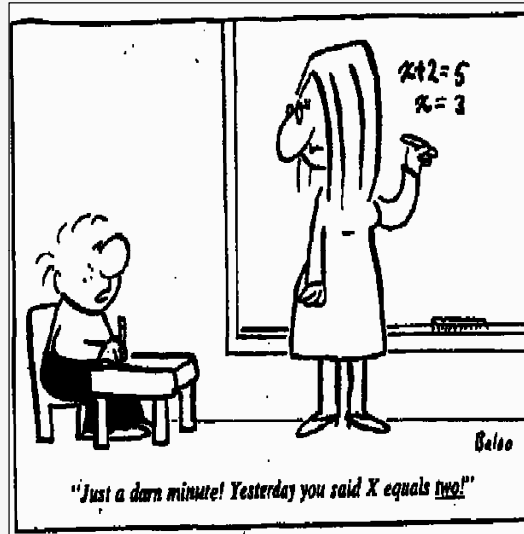
Table 1. Descriptive statistics (N = 1131)

Category and variable	Mean	Range	Standard deviation
<i>Background</i>			
Sex (1 = female)	0.488	0-1	0.499
Social class	1.948	1-3	0.812
Residence (1 = Reykjavik area)	0.444	0-1	0.494
<i>Attitudes/beliefs</i>			
Importance of sport	0.863	0-2	0.744
Importance of health improvement	1.805	0-2	0.445
Internal health locus of control	3.437	0-4	0.825
<i>Other's physical activity</i>			
Father	0.732	0-2	0.822
Mother	0.740	0-2	0.830
Older brother	0.603	0-2	0.767
Older sister	0.462	0-2	0.688
Best friend	1.448	0-2	0.755
Main teacher	1.352	0-2	0.677

Constant vs. Variable

Variables are the specific properties that have the ability to take different values.

Constants are the specific properties that cannot vary or won't be made to vary.



9

Terminology - Variables

INDEPENDENT

(syn: treatment, experimental, predictor, input, exposure, explanatory variable) is a stimulus or activity that is *identified or manipulated* to predict the dependent variable; they are considered as the causal factors, or that you may manipulate.

e.g. new drug, working hours, exposure, worker attitudes, policies

DEPENDENT

(syn: Effect, criterion, criterion measure, outcome, output variable) is a response that the researcher wanted to predict; they are considered as the *outcomes of the treatments or the responses* to changes in the independent variables.

e.g. Symptomatology, productivity, accident rates, attitudes, health status, performance on neuropsychological test

Adapted from: Dr. Craig Jackson, University of Central England

UCE
Birmingham

Terminology - Variables

EXTRANEIOUS

Extraneous variable is a variable that has a potential to *distorts the relationship between dependent and independent variables*.

- **Controlled extraneous variables** are recognized before the study is initiated and are *controlled in the design and selection criteria*.

- **Uncontrolled extraneous variables** are recognized before the study is initiated or, sometimes, even if recognized cannot be controlled in the design and selection phase. Usually an attempt is made to *assess and adjust them through sophisticated statistical tools*.

e.g., Working hours, temperatures, extraneous exposure, diet, class, income, Ambient noise and temperature in testing room

Adapted from: Dr. Craig Jackson, University of Central England

UCE
Birmingham

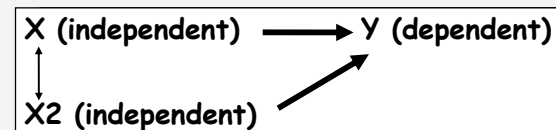
Study Variables

Independent Variables & Dependent variables

$X \longleftrightarrow Y$

$X \text{ (independent)} \longrightarrow Y \text{ (dependent)}$

Extraneous variable



12

Study Variables

Confounding Variable

- When the *effects of two or more variables cannot be separated.*

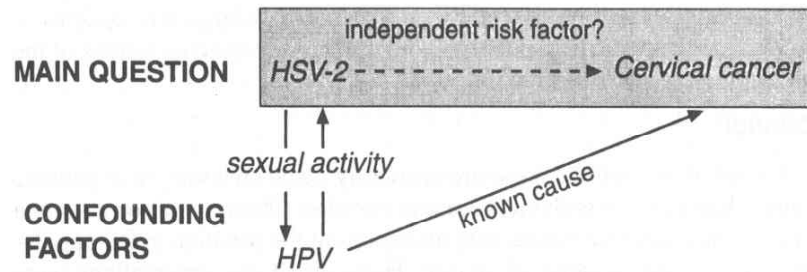


Figure 1.1. Confounding bias: Is herpesvirus 2 (HSV-2) a possible cause of cervical cancer? Only if its association with cervical cancer is independent of human papillomavirus (HPV) infection, known to be a cause of cervical cancer. Both viruses are related to increased sexual activity.

13

Study Variables

Confounding Variable

- When the *effects of two or more variables cannot be separated.*

		STD rate	
Condom Use	Yes	55/95	(61%)
	No	45/105	(43%)

"Condom Use increases the risk of STD"

BUT ...

		STD rate	
# Partners < 5			
Condom Use	Yes	5/15	(33%)
	No	30/82	(37%)
# Partners > 5			
Condom Use	Yes	50/80	(62%)
	No	15/23	(65%)

Explanation: Individuals with more partners are more likely to use condoms. But individuals with more partners are also more likely to get STD.

14

Bias & Chance

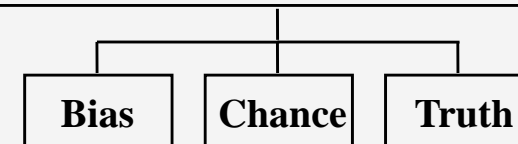


15

Measuring Outcomes:

Observed vs. Truth

Possible Explanations of Outcome Measured



$$\text{Observed} = \text{Truth} + \text{Error}$$



Systematic error + Random error
(Bias) (Chance)

16

Bias vs. Chance

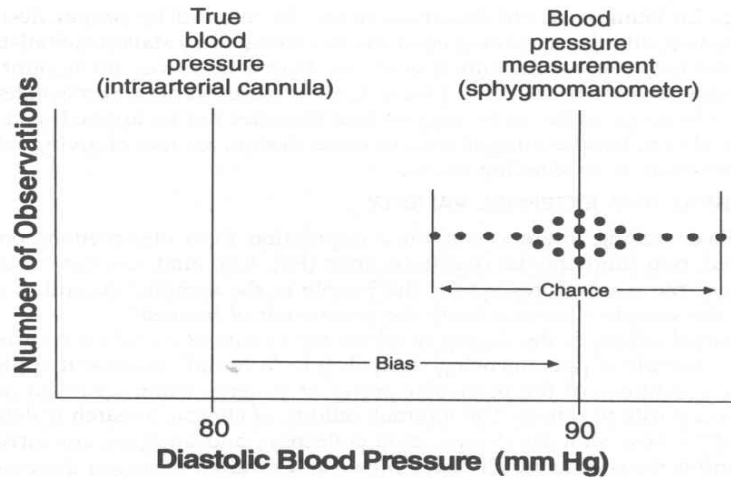


Figure 1.2. Relationship between bias and chance: Blood pressure measurements by intraarterial cannula and sphygmomanometer.

Bias:

- A process at any stage of inference tending to produce results that depart systematically from the true values.

Chance:

- The divergence of an observation on a sample from the true population value in either direction.
- The divergence due to chance alone is called *random variation*



Bias and chance- are not mutually exclusive. 18



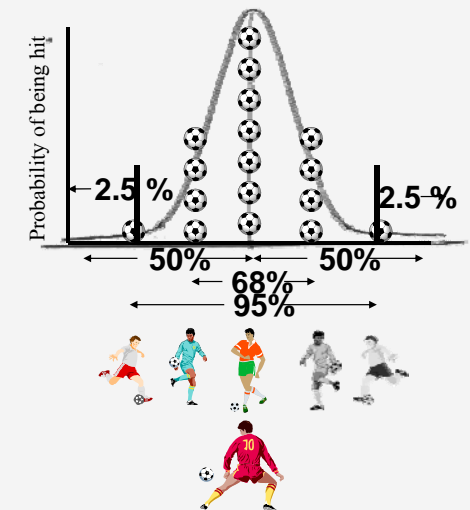
“A well designed, carefully executed study usually gives results that are obvious without a formal analysis and if there are substantial flaws in design or execution a formal analysis will not help.”

Johnson AF. Beneath the technological fix. *J Chron Dis* 1985 (38), 957-964 19

Chance



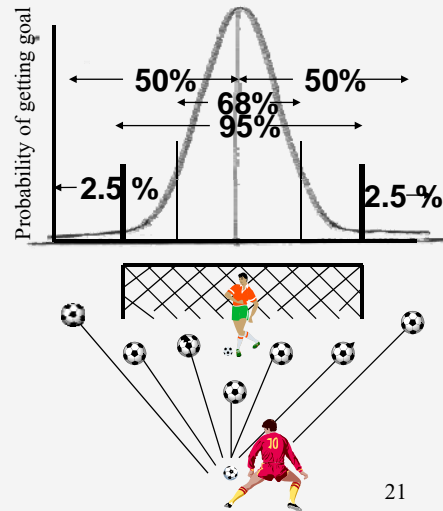
“Free kick”



Chance



"Free kick"



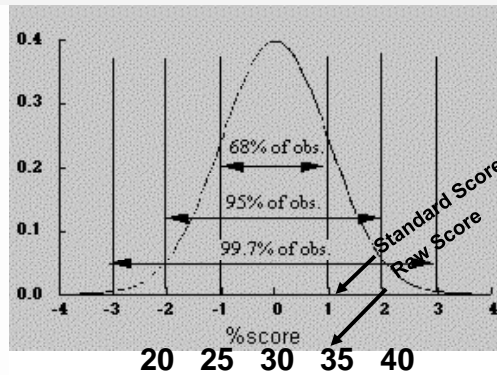
21

Normal Distribution



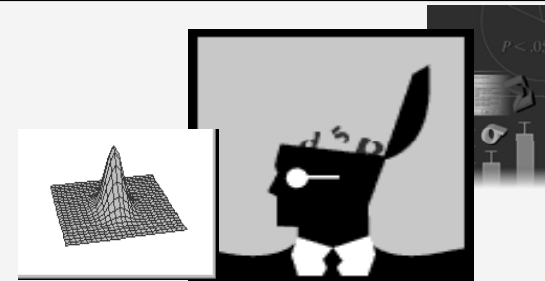
22

Normal Distribution in Descriptive Statistics



$X = 30;$
 $SD = 5$

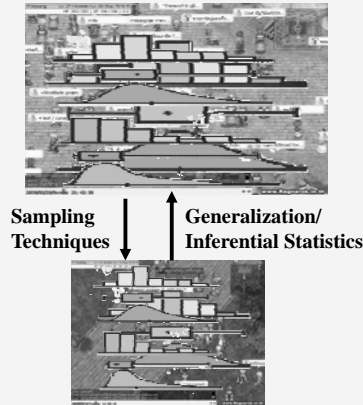
Types of Statistical Methods



24

Types of Statistics

- By Level of Generalization
 - Descriptive Statistics
 - Inferential Statistics
 - Parameter Estimation
 - Hypothesis Testing
 - Comparison
 - Association
 - Multivariable data analysis
- By Level of Underlying Distribution
 - Parametric Statistics
 - Non-parametric Statistics



25

Descriptive Statistics



26

Descriptive Statistics

• Measure of Location (Categorical Vars)

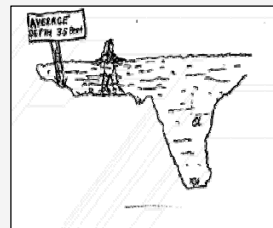
- Frequency (f)



• Measure of Location (Continuous Vars)

- Mean
Average

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \quad \text{or} \quad \mu = \frac{\sum_{i=1}^n x_i}{N}$$



- Median
Mid-point

$$x_1 x_2 x_3 x_4 x_5 x_6 x_7 x_8 x_9$$

- Mode

The Most Frequent

$$x_1 x_2 x_3 x_4 x_5 x_6 x_7 x_8 x_9$$

(1 2 2 2 2 3 3 4 5)

27

Descriptive Statistics

• Measure of Spread

- Range

Max - Min

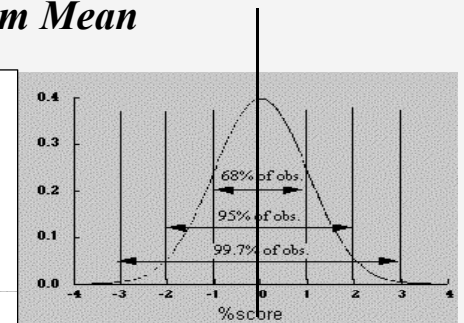
- Standard Deviation / Variance

X_i s deviate from Mean

$$S = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n - 1}}$$

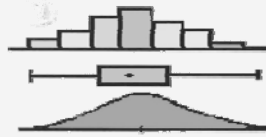
or

$$\sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \mu)^2}{N}}$$



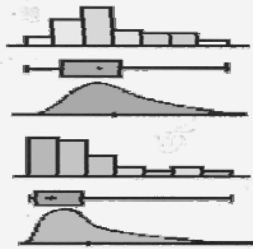
- **Measure of Shape**

- Normal Distribution

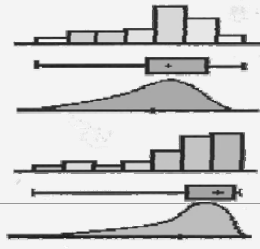


- Skewed Distribution

- Positively skewed



- Negatively skewed



Inferential Statistics



30

- **Purpose of Inferential Statistics**

- Generalisability of research results from Sample Statistics to Population Parameters

- **Types of Inferential Statistical Methods**

- **Parameter Estimation** - to estimate the range of values that is likely to include the true value in population

$$\bar{X} \rightarrow \mu$$

$$\text{proportion} \rightarrow \pi$$

- **Hypothesis Testing** - to ask whether an effect (difference) is present or not among different groups

$$H_0: \bar{X}_1 = \bar{X}_2 \rightarrow H_0: \mu_1 = \mu_2$$

$$H_0: r_{xy} = 0 \rightarrow H_0: \rho_{xy} = 0$$

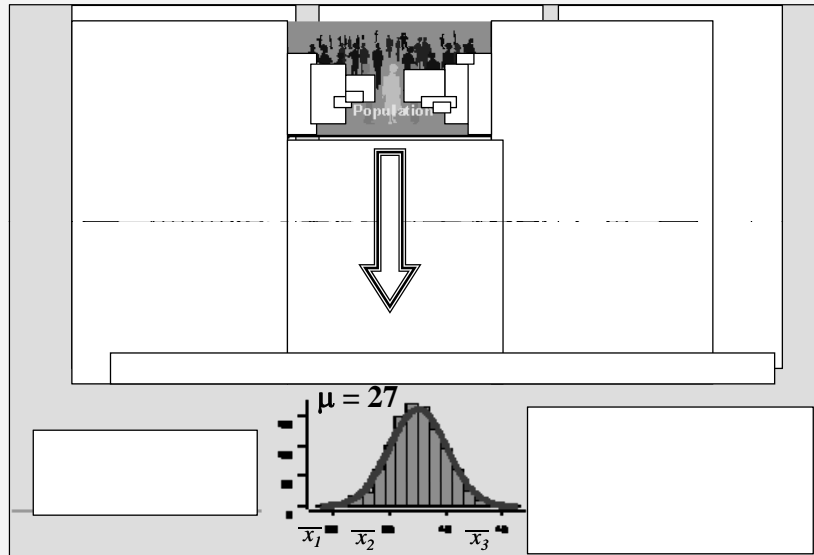
Herman



"How could I have been doing 70 miles an hour when I've only been driving for ten minutes?"

32

Sampling Distribution



Sampling Distribution (Normal Distribution) in Parameter (μ) Estimation

Confidence Limits of μ : $\bar{X} \pm Z_{\alpha/2, v} S_{\bar{X}}$

Standard error
95% CI of μ : $\bar{X} \pm (1.96 * (SD/\sqrt{n}))$
 $25 \pm (1.96 * (12.2/\sqrt{100}))$
 $\mu = 25$ (22.6 to 27.4, 95% CI)

Average = 26
SD = 9.1

Average = 30
SD = 11.3

Average = 25
SD = 12.2

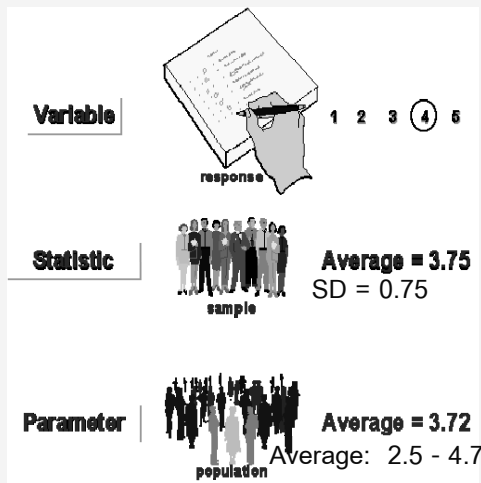
The Sampling Distribution...

$\mu = 27$ $N = 10,000$

...is the distribution of a statistic across an infinite number of samples

Point Estimates:

Single values (Mean, Variance, Correlation, treatment effect, relative risk, etc.) representing characteristics in the whole population



Interval Estimates:

Ranges of values, usually centered around point estimates, indicating bounds within which we expect the true values for the whole population to lie (stability of the estimate)

Example Confidence Limits of μ : $\bar{X} \pm Z_{\alpha/2, v} S_{\bar{X}}$ 35

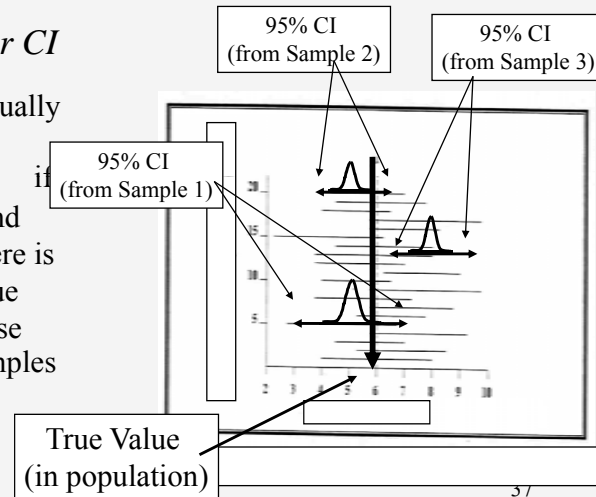
Point estimates and confidence intervals are used to characterise the statistical precision of any rate (incidence, prevalence), comparisons of rates (e.g., relative risk), and other statistics.

- US adults have used unconventional therapy = 34% (31% - 37%, 95%CI)
- Sensitivity of clinical examination of splenomegaly = 27% (19 - 36%, 95%CI)
- Relative risk of lung cancer of smoker vs. non-smoker = 5.6 (2.1 - 8.9, 95%CI)
- Relative risk of HIV infected of male vs. female = 2.1 (0.5 - 6.9, 95%CI)

Consideration in confidence level of the estimate

Select a cut-point for CI

Confidence Interval, usually set at 95% CI, can be interpreted such that - if the study is unbiased and repeated 100 times, there is 95% chance that the true value is included in these intervals of the 100 samples



37



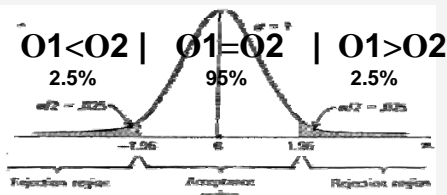
38

- **Hypothesis & Tail of the test**
– One-sided vs. Two-sided Test

Two-sided test:

Ex: H_0 : Outcome 1 = Outcome 2

H_a : Outcome 1 \neq Outcome 2



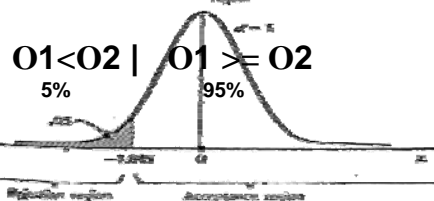
One-sided test:

Ex: H_0 : Outcome 1 \leq Outcome 2

H_a : Outcome 1 $>$ Outcome 2

H_0 : Outcome 1 \geq Outcome 2

H_a : Outcome 1 $<$ Outcome 2



39

- Basic steps in hypothesis testing

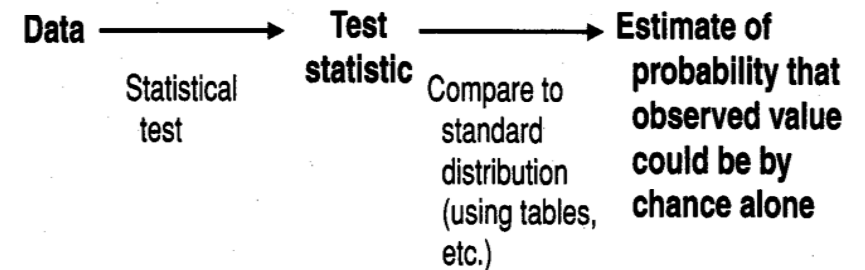
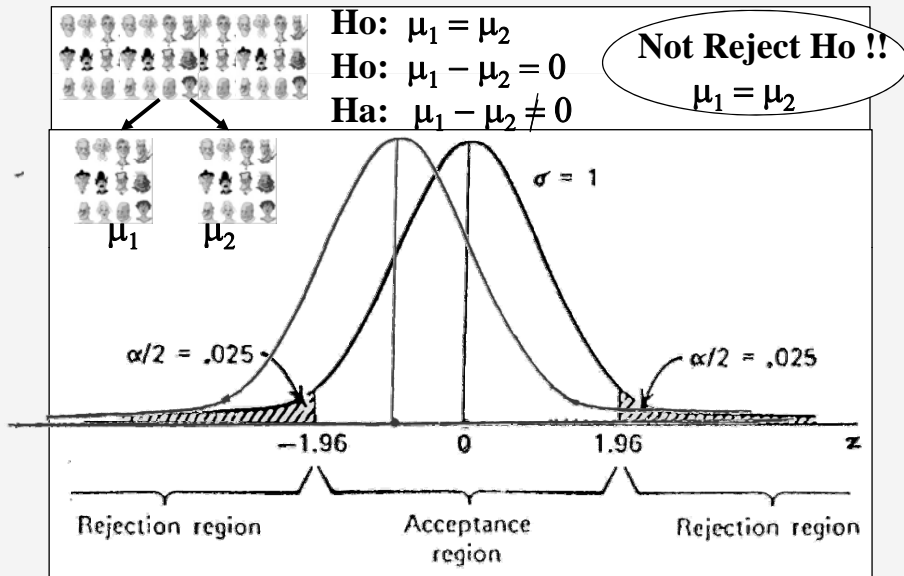


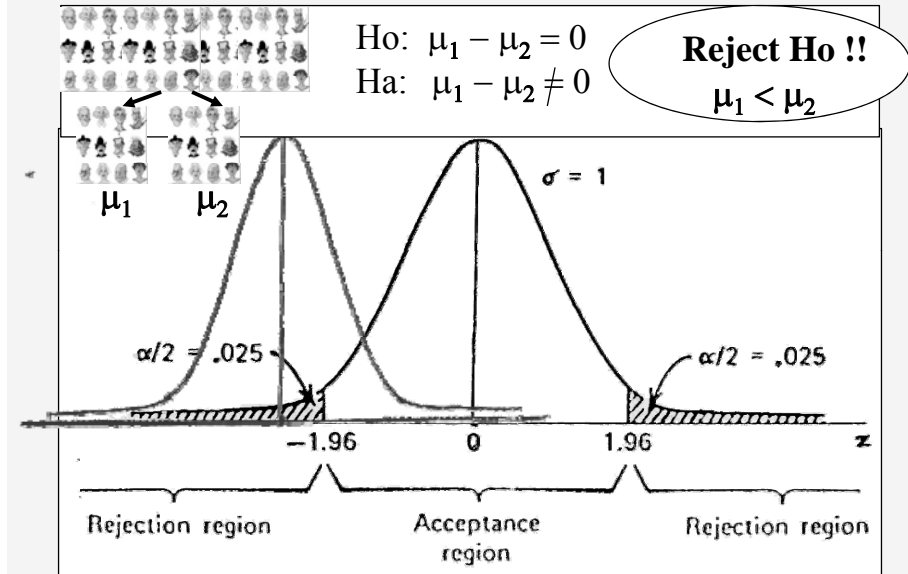
Figure 9.2. Statistical testing.

40

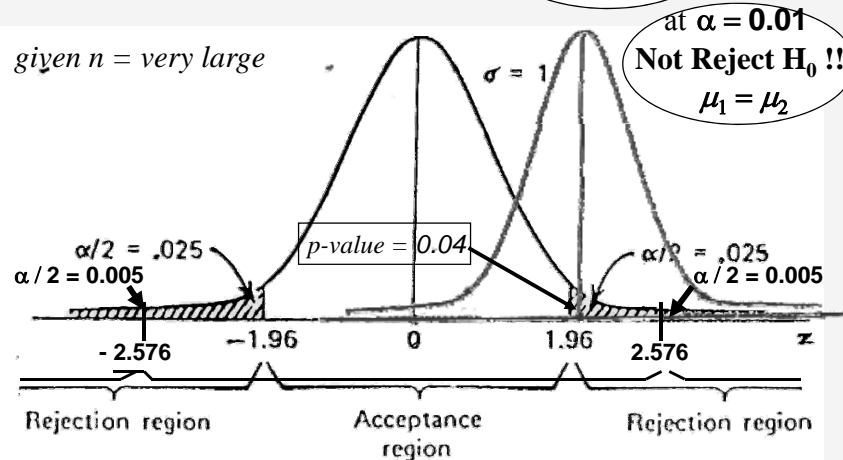
Hypothesis Testing



Hypothesis Testing



$H_0: \mu_1 - \mu_2 = 0$
 $H_a: \mu_1 - \mu_2 \neq 0$



• Type I & Type II errors

$H_0: G_1 = G_2$

Reality/Truth

		Reality/Truth	
		H_0 True ($G_1 = G_2$)	H_0 False ($G_1 > G_2$)
Decision	Accept H_0 ($G_1 = G_2$)	Correct Confidence : $1 - \alpha$ 0.99, 0.95	Type II Error β 0.10, 0.20
	Reject H_0 ($G_1 > G_2$)	Type I Error α 0.01, 0.05	Correct Power : $1 - \beta$ 0.90, 0.80

The End of Inferential Statistics

