Week 1

Monday, 2 August 2021 23:36

BIOSTATISTICS

IMPORTANT FOR PUBLIC HEALTH EDUCATION, MESEANCH & PRACTICE

STUDYING SUB-SET OF A LANGER POPULATION OR PROCESS

STUDY DESIGN PRACTICES TYPES OF DATA WE WILL ENCOUNTER

STATISTICS
DATA SCIENTISTS
BIG DATA ENGINEER

SEXIEST SOBS
OF THE 21TH C.

DATA IS EVERYWHERE DATA IS THE NEW OIL

PLANNING
DESIGN OF STUDY
DATA COLLECTION
DATA ANALYSIS
PRESENTATION
INTERPRETATION

DATA
PIPELLNE
RESEARCH
MOCESS

· SAMPLES US POPULATION

IMPENPECT NAPPNESENTATION OF SOME LANGER POPULATIONS

SUBSET OF POPULATION FROM WHICH INPO IS COLLECTED

II IMPENFECT NEAU ZATION"

WE CANNOT OBSEING POPULATION MINECTLY

ENNOR ASSOCIATED WITH CHOOSING THE SAMPLE

SAMPLE SHOULD BE NEPHESENTATIVE

SIMPLE NANDOM SAMPLING

AVENY POSSIBLE SUBSET OF SIZE N FROM PUPULATION IS EQUALLY LIKELY TO LE SELECTED IF SAMPLE NOT NEPHRESENTATIVE --- BIAS

SOMETIME THIS IS NECESSARY

EX DAVO USER IN MUMBAI INDIA (LIST ON NEGISTRY)

WOMEN AME LESS LIKELY TO BE OPEN ALOUT THEIR MING USE

RANDOM DIGIT DIALLING COMEMENCE SAMPLING NESPONDENT DINVEN SAMPLING USED WHEN
SAMPLE IS
DIPPICULT TO
OBTAIN

NOT ALL ELEMENTS OF A STUDY CAN DE SAMPLED

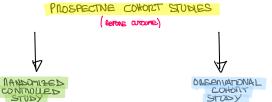
STATISTICS

HON TO DEAL WITH THE IMPENFECT OF NEW APPLIATION

HOW CAN WE ACCOUNT FOR THE UNCENTAINTY IN OUR ESTIMATE?

SAMPLING PROCESURE?





NANDOMLY ASSIGN NAMPLE MEMBERS TO EX POSUNE GNOVPS

SMONE CAUSES XII

ONLY SYSTEMATIC DIPFENENCE IS EXPOSED OF INTENEST ONOUP MEMBERSHIP, MEASURING EPFECT ON APFERENT CROUPS

I SMOKERS VS NOW SMOKERS 211
PERFORM DIFFERENCY

SELF-SELECTION

METHOD MAY LEADS TO SYSTEMATIC DIFFERENCES

(I MAYBE SMOKENS ARE MORE LIKELY TO DRINK QUONTUL AND BOTH ARE LUMED TO CANDOVISCUAL MEDADERS

المستحدية ما

MINIECHON AGAINST SELF SELECTION BLASES

POTENTIAL

THE OUTCOME IT EXPOSITE

THEATMENT DECIDED BY NESEANCHES

COMPLEX NANDOMIZINO SCHEMA

(NOT ALWAYS POSSIBLE)

CASE/CONTROL STUDIES (AFTER OUTCOME)

IN PRECEDENT STODIES SOLDECT WERE ASSIGNED AN EXPOSURE STATUS ON WERE SELECTED A THEIR EXPOSED STATUS WAS CLASSIPLED

EXPOSURE W/ NAME OUTCOMES

SELECTED PEOPLE ON WHETHER THEY HAVE AN OUTCOME ON NOT (MOSTLY A NAME DISEASE)

ANALYTICAL ISSUE

CAN'T ESTIMATE NISK OF OUTCOME

INPURENCE BY HOW NESEANCHES CHOOSE CONTÁGUS & CASES

EX LINK LETWEEN SMOKE IT WING-CANCER

POTENTIAL CONFOUNDERS

THERE CAN BE SYSTEMATIC DIFFERENCES BESIDE EXPOSURE DIFFERENCES

WHAT CONFOUNDERS HAVE NOT BEING ADDRESSED?

NECALL BIAS

PATIENT INTERVIEW (THEY TEND TO EXAGENATE)

· DATA TYPES & SUMMEDIZATION

CONTINUOUS DATA (INCREMENTAL) MMMg.H. [NO.], Elyran

Yes-No

CATEGORICAL DATA (NOMINAL & ONDINAL)

NACE, GENDER, NEUGION (NOMINAL)

ONDER (STRONGLY) NO 01/10/2012

BINARY IS A SPECIAL CASE OF CATEGORICAL DATA (2 LEVELS)

TIME-TO-EVENT DATA

#OFTIMES, PLACE, TIME

COMPANING GLOOD PRESSURE?

- ESTIMATE M NEAN DIPPEMENCE IN LA CF GE IN EACH OF THE GROUP
- ESTMATE A 95% CONPIDENCE INTERNAL & USE A T-TEST FOR POPULATION LOUGL DIFFERENCES
 - · CONFIDENCE INTERNALS

 - · PVALUE · T, 2, 12, F TESTS · KAPLAN MOLEIL

 - · LOG-NANK TEST
 - · PROPORTIONS, PROP TEST

DIPFENENT APPROACHES TO ANALYZE DIFFERENT DATA-TYPES

DATATYPES IT STUDY-DESIGNS