

## Sampling & Designing

→ How is data produced, drawing samples & designing experiments

### Statistical Inference :-

Population: the entire population / group of subjects about which we want information (e.g.: total Indian population)

Parameter : quantity about population we are interested in  
(e.g.: approval % among all Indian voters)

Sample : part of population from which we collect information  
(e.g.: 1,000,000 sample size)

Statistic / Estimate: the quantity we are interested in as measured in sample  
(approval % among the sample voters).

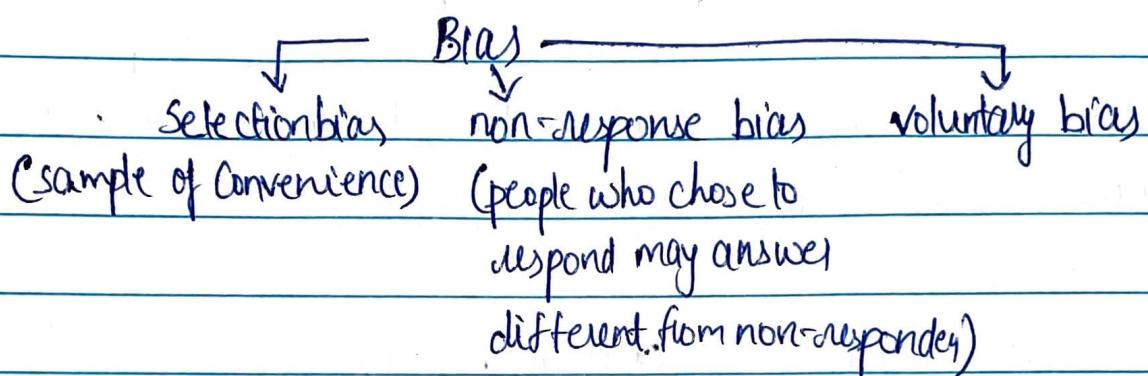
⇒ Even a relatively small sample will produce an estimate that is close to the parameter of a very large population.

## Simple random sampling & other sampling plans:-

→ Sampling Correctly is very important.

→ Selecting 1000 from your neighbour, is sample of Convenience.

This is not a good way as it will introduce bias, sample favoring certain outcome



Ways to avoid these bias is to use chance in a planned way.

→ A simple random sample, selects sample at random without replacement.

→ A stratified random sample, divides population into groups of similar subjects called "strata" (e.g. rural, urban). Then a simple random sample is chosen. This is resulting in a more precise estimate.

→ Since sample is drawn at random, estimate will be different from the parameter due to "chance error". Another sample gives another "chance error". By taking a bigger sample size, "chance error" can be made small.

$$\text{estimate} = \text{bias} + \text{chance error} + \text{parameter}$$

But if we increase the sample size, the bias will not get smaller & typically we don't know how large bias is.

## Randomized controlled experiments & observation studies:-

Let's say "people who eat red meat has a chance of getting certain cancers than who don't".

- This is called "association".
- But this does not mean, eating red meat causes cancer.
- Association doesn't mean causation.
- This is an observational study. It measures outcomes of interest & this can be used to establish association
- Because there maybe "Confounding factors" or "lurking variable"
- To establish causation, an experiment is required.

A treatment (eating red meat) is assigned to people in treatment group and other is control group. To the end, the confounding factors should be same, then compare outcomes.

- ① The best way to make sure that 2 groups are similar is to assign a subject at random.
- ② Control group gets placebo.
- ③ Experiment is double-blind. Neither subject nor evaluator know which subject is in which group.

⇒ A good experiment requires that the subjects are assigned groups at random so that they operate equally, apart from differences due to chance. Since it is chance effect, we can calculate the size when we evaluate the outcome.

### Interpretation of Probability

Probability of an event is defined as proportion of times this event occurs in many repetitions, given it is possible to repeat.

e.g. In 2015, 4 million babies were born & 48.8% were female.

$$\Rightarrow P(\text{newborn is female}) = 48.8\%.$$

→ Long interpretation of probability can make it difficult to interpret it for single event.

→ Subjective probability is not based on experiments but on subjective probabilities of an event by individuals.

Key for Computing probabilities:

① Probabilities are always between 0 & 1.

② Complement rule  $\Rightarrow P(\bar{A}) = 1 - P(A)$

③ Rule for equally likely outcome ( $n$ )  $\Rightarrow$

$$P(A) = \frac{\text{no. of outcomes in } A}{n}$$

④ A & B are mutual exclusive if they cannot occur same time

Addition rule  $\Rightarrow P(A \text{ or } B) = P(A) + P(B)$

multiplication  $\Rightarrow P(A \& B) = P(A) P(B)$  for this two

events should be independent, i.e. knowing one occurs does not change the probability of others.

⑤ P(at least one 6 in 3 rolls).

$\Rightarrow$  write 'at least one 6'

first roll or second roll or third roll

but these are not mutually exclusive, as we can't have 6 on second roll.

$$\therefore P(\text{at least one 6}) = 1 - P(\text{no 6 in 3 rolls})$$

$$= 1 - P(\text{no 6 in first and second and third})$$

$$= 1 - \frac{5}{6} \times \frac{5}{6} \times \frac{5}{6}$$

$$= \underline{41.4\%}$$

## Conditional probability :-

→ A spam email contains the word 'money' higher chance than a ham/normal email.

$$P(\text{'money' in email} \mid \text{spam}) = 8\%. \quad " \mid " \rightarrow \text{given that}$$

$$P(\text{'money'} \mid \text{ham}) = 1\%.$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$\Rightarrow \frac{P(B \cap A)}{P(A \cap B)} = P(A) P(B|A) = P(B) P(A|B)$$

if A & B are independent,  $P(A \cap B) = P(A) P(B) = P(B) P(A)$

① Let  $P(\text{spam}) = 20\%$ . what is the probability if contains 'money' in an email.

$$\textcircled{1} \quad P(\text{spam}) = 20\%, \quad \Rightarrow P(\text{ham}) = 80\%.$$

$$P(\text{money} \mid \text{spam}) = 8\%.$$

$$P(\text{money} \mid \text{ham}) = 1\%.$$

Event that "money" appears in email can be written as

money and email is spam (or) | money and email is ham

⇒ these 2 events are mutually exclusive, i.e if email is spam, then it can't be ham

$$\therefore P(\text{money appears}) = P(\text{money} \& \text{spam}) + P(\text{money} \& \text{ham})$$

$$\Rightarrow (8 \times 20) + (1 \times 80) = P(\text{money} \mid \text{spam}) P(\text{spam}) + P(\text{money} \mid \text{ham}) P(\text{ham})$$

$$\Rightarrow 160 + 80 \Rightarrow \frac{240}{100} \Rightarrow \underline{\underline{24\%}}$$

## Bayes' Rule :-

From data we know  $P(\text{money}|\text{spam}) = 8\%$ , but what we need to build is  $P(\text{spam}|\text{money})$

$$\text{(Bayes' Rule)} \quad P(B|A) = \frac{P(A \text{ and } B)}{P(A)} = \frac{P(B \text{ and } A)}{P(A)} = \frac{P(B) P(A|B)}{P(A)}$$

$$\Rightarrow P(\text{spam}|\text{money}) = \frac{P(\text{money}|\text{spam}) P(\text{spam})}{P(\text{money})} = \frac{8 \times 20}{240} = \frac{2}{3} = \underline{\underline{67\%}}$$

$$P(B|A) = \frac{P(A \text{ and } B)}{P(A)} = \frac{P(B \text{ and } A)}{P(A)} = \frac{P(B) P(A|B)}{P(A)}$$

$$= \frac{P(A|B) P(B)}{P(A|B) P(B) + P(A|\bar{B}) P(\bar{B})}$$

## Bayesian Analysis:-

The spam filter classifies email as spam via Bayesian Analysis

$$\rightarrow P(\text{spam}) = 20\%$$

$\rightarrow$  After examining keywords like 'money', the posterior probability is improved.

⑥ 1% population has disease, test shows 95% positivity rate if infected is tested. If non-infected is tested, the test shows 2% chance as positive. Given person has tested positive, what is probability they are infected?

$$\textcircled{1} \quad P(\text{disease is present}) = 1\% \Rightarrow P(\text{disease not present}) = 99\%.$$

$$P(\text{Positive} | \text{disease is present}) = 95\%,$$

$$P(\text{positive} | \text{disease is not present}) = 2\%.$$

$$P(\text{disease is present} | \text{positive}) = ?$$

$$\therefore P(\text{disease} | \text{positive}) = \frac{P(\text{positive} | \text{disease}) P(\text{disease})}{P(\text{positive})}.$$

$$\Rightarrow P(\text{positive}) = P(\text{positive} | \text{disease}) P(\text{disease}) + P(\text{positive} | \text{no disease}) P(\text{no disease})$$

$$= (0.95 \times 0.01) + (0.02 \times 0.99)$$

$$= 0.0293$$

$$\therefore P(\text{disease} | \text{positive}) = \frac{0.95 \times 0.01}{0.0293} = 0.0324 = \underline{\underline{32.4\%}}$$

Reason is, only 1% of population has disease & test finds most of those. But there's 2% chance of not infected. Since that population is higher, the probability is low.

## Warner's Randomized response model :-

Assume surveying "what % students have cheated during online exam".

But as students won't answer truthfully, we use Randomization as:

$$\rightarrow (HH, HT, TH, TT)$$

Toss a coin twice, if ~~if~~ they get "T", then they answer question-1,  
else question-2.

Q-1 → did you cheat in exams

Q-2 → did you get tails on 2<sup>nd</sup> toss

So the answer will be partly random. A 'yes' could be to any of these questions.

→ While we don't know what individual 'yes' means, we can estimate the proportion of cheaters using all answers collectively.

$$\begin{aligned} P(\text{yes}) &= P(\text{yes} \cap Q_1) + P(\text{yes} \cap Q_2) \\ &= P(\text{yes}|Q_1) P(Q_1) + P(\text{yes}|Q_2) P(Q_2). \end{aligned}$$

$$\Rightarrow P(\text{yes}|Q_1) = \frac{P(\text{yes}) - P(\text{yes}|Q_2) P(Q_2)}{P(Q_1)}$$

Let 27 answered yes & 30 said no, then

$$P(\text{yes}) = \frac{27}{27+30} = 47\%.$$

$P(Q_2) \rightarrow$  just half  
because it depends  
on first coin toss

$P(Q_1) \rightarrow$  half

$P(\text{yes}|Q_2) \rightarrow$  half

$$P(\text{yes}|Q_1) = \frac{0.47 - (0.5 \times 0.5)}{0.5} = 0.44 = \underline{\underline{44\%}}$$

Q) A fair coin is tossed 5 times.  $P(\text{at most 4 tails})$ .

A)  $P(\text{at most 4 tails})$  can be written as

$P(1 \text{ tail}) \quad P(2 \text{ tail}) \quad P(3 \text{ tail}) \quad P(4 \text{ tail})$  (not mutually exclusive)

$$\Rightarrow 1 - P(\text{no tail in 4 tails})$$

$\Rightarrow 1 - P(\text{no tail in 1st tail}) \text{ and } P(\text{2nd tail}) \text{ and } P(\text{3rd tail}) \text{ and } P(\text{4th tail})$

$$\Rightarrow 1 - \left(\frac{1}{2}\right)^5$$

Q) 3 boxes  $\rightarrow$  Box 1  $\rightarrow$  2 quarters

Box 2  $\rightarrow$  2 nickels

Box 3  $\rightarrow$  1 quarter & 1 nickel.

Pick a random  
Coin from a  
random box.

If coin picked is quarter, what is chance other coin is also quarter in  
box.

$$A) P(\text{selecting Box 1}) = P(\text{selecting Box 2}) = P(\text{selecting Box 3}) = \frac{1}{3}$$

$$\Rightarrow P(\text{quarter} | B_1) = 1 \quad P(\text{quarter} | B_2) = 0 \quad P(\text{quarter} | B_3) = \frac{1}{2}$$

$$\Rightarrow P(B_1 | \text{quarter}) = \frac{P(\text{quarter} | B_1) P(B_1)}{P(\text{quarter})}$$

$$P(\text{quarter}) = P(B_1) P(\text{quarter} | B_1) + P(B_2) P(\text{quarter} | B_2) + P(B_3) P(\text{quarter} | B_3)$$

$$= \left(\frac{1}{3} \times 1\right) + \left(\frac{1}{3} \times 0\right) + \left(\frac{1}{3} \times \frac{1}{2}\right) = \frac{1}{3} + \frac{1}{6}$$

$$\Rightarrow P(B_1 | \text{quarter}) = \frac{\left(\frac{1}{3} \times \frac{1}{2}\right)}{\left(\frac{1}{3} \times 1\right) + \left(\frac{1}{3} \times \frac{1}{2}\right)}$$

Q) Roll pair of dice. What is the chance of double when you roll a pair of dice?  $(1,1) (2,2) (3,3) (4,4) (5,5) (6,6)$

A) mutually exclusive and independent.

$P(\text{Showing double})$

$$\Rightarrow 1 - P(\text{not showing double})$$

$\begin{array}{c} (1,2) (1,3) \dots (1,6) \rightarrow 5 \\ (2,1) \\ (3,1) \\ (4,1) \\ (5,1) \\ (6,1) \end{array} \quad \Rightarrow 30$

$$\Rightarrow 1 - \frac{5}{30} \Rightarrow \frac{1}{6} \quad (\text{Out of 30 possible outcomes, 6 are double})$$

Q) Monopoly, if you land in jail, must roll a double on any of next 3 turns, or pay fine. - what is  $P(\text{get out of jail} | \text{no fine})$ ?

A)  $P(\text{double in next 3 turns})$

$P(\text{in 1st turn}) \text{ or } P(\text{in 2nd turn}) \text{ or } P(\text{3rd turn})$  (not mutually exclusive)

$$\Rightarrow \frac{1}{6} \text{ or } \frac{1}{6} \text{ or } \frac{1}{6} \quad \text{Probability of rolling}$$

$$\Rightarrow 1 - \left(\frac{5}{6}\right)^3 \left(\frac{5}{6}\right)^3 \quad \text{at least 1 double is same as not rolling 3 non-doubles}$$

Q) 3% of applicants are admitted. 70% of all applicants have GPA  $3.6 \leq$  of those admitted, 95% have GPA  $3.6 \leq$ . Chance of applicant being admitted with GPA  $\geq 3.6$ ?

$$P(\text{admit} | 3.6) = \frac{P(\text{admit} \cap 3.6)}{P(3.6)}$$

A)  $P(\text{admitted}) = 3\% \Rightarrow P(\text{not admitted}) = 97\%$

$$P(3.6 | \text{admitted}) = 70\% = P(\text{applicant} | 3.6)$$

$$= \frac{P(3.6 | \text{admitted}) P(\text{admitted})}{P(3.6)}$$

$P(3.6 | \text{admitted}) = P(\text{admitted} | 3.6) = 95\%$ , then  $P(3.6 | \text{admitted}) = 95\%$ .

② MCQ has 10 questions. Each question has 3 possible answers, one is correct. A student knows correct answer to 4 & guesses other. If then first question was correct, what is chance of it being guess.

$$\text{③ } P(\text{answer is correct}) = \frac{1}{3}$$

$$P(\text{answer is correct} | \text{know}) = \frac{4}{10}$$

$$P(\text{answer is correct} | \text{guess}) = \frac{6}{10}$$

$$\Rightarrow P(\text{guess} | \text{answer}) = \frac{P(\text{answer} \& \text{guess})}{P(\text{answer})}$$

$$= \frac{P(\text{know}) P(\text{answer} | \text{guess}) P(\text{guess})}{P(\text{answer}) P(\text{answer} | \text{guess}) + P(\text{guess}) P(\text{answer} | \text{know})}$$

$$= \frac{\frac{1}{3} \times \frac{6}{10}}{(1 \times \frac{6}{10}) + (1 \times \frac{4}{10})}$$

$$\underline{\underline{\frac{\frac{1}{3} \times \frac{6}{10}}{(1 \times \frac{6}{10}) + (1 \times \frac{4}{10})}}}$$