

- **Anecdotal Evidence : Evidence based from a limited sample size which might not be representative of the population**

First question you should ask : What is the population and what is the sample?

Research Question \rightarrow Population \rightarrow Sample \rightarrow Generalize to

Type of variable:

- Numerical: Takes numerical values (Add, subtract on these)
 - Continuous : E.g height
 - Discrete: E.g. age
- Categorical : Limited number of distinct categories, can be identified with numbers but no arithmetic operation
 - Ordinal: Have ordered levels , Satisfactory, poor, very poor
 - Regular : Morning person or Afternoon person

When variables show some connection with one another they are called associated or dependent can be neg or positive . Not associated = independent

- Observational Study: Collect data without affecting how data arises, "observe"
 - Retrospective: Data collected from past
 - Prospective: Collected through study
- Experimental Study: Randomly assign subjects to treatment, thus can establish causal connections between the explanatory and observed variables

In an observational study it is difficult to conclude, unless you really control the effect of other variables. In experimental study due to random assignment this is taken care of.

Confounding variables: Variables that affect both explanatory and response variables

Why is 'census' not a good idea?:

- Part of the 'census' might not be representative of the population, e.g. Immigrants
- Population is dynamic! You don't taste whole of soup to find out its taste! Extrapolatory analysis!
- For inference to be valid, the sample should represent population. Stir the soup before tasting!

Bias:

- Convenient Sample: Pick up people from your class for study because they are easily accessible
- Non response: Non-random section of people respond to your survey. Emailing a survey to people who do not have internet connection \Rightarrow No point!

- Voluntary Response: Only people who volunteer to respond, respond. not everyone

The Literary Digest shut down because its sample was not representative of the population, however large the sample was! Sampling :

- Simple Random Sampling: Randomly select samples so that each data point has equal probability to get sampled
- Stratified: Divide population into homogenous groups(strata) and sample from these. Eg. if you want males and females to be equally represented divide them into males and females and then sample
- Cluster: Divide into non-homogenous cluster, sample the clusters and then sample the data points. For eg.e divide the geography into clusters , travel to few clusters only

Principles of experiment design:

- Control: Compare treatment of interest to group
- Randomize: random assignment of subject to treatment
- Replicate; replicate entire study or more samples
- block: Block variables that might affect response variables. Divide pro and amateur athletes into two groups, then assign them independently to treatment and control and then observe.

Blocking vs explanatory

- Explanatory variables/factors are conditions we can impose on the experimental units
- Blocking variables are characteristics that experimental units come with that we would like to control
- placebo: fake treatment
- placebo effect: experiments show effect just because they think they are undergoing treatment
- blinding: subjects don't know whether they are in control or sample
- double blind: neither the experimentalist nor subject knows the groups

Random sampling" Randomly select subjects from sample Random assignment: randomly assign subjects to treatment and control, thus removing the chance of difference

First you sample randomly from the population and then randomly assign them to control and treatment groups

RS, RA => Causal generalisable [IDEAL experiment] NO RS, RA => Causal not generalizable [Most Experiments] RS, NO RA => Not causal but generalizable [MOST Observations] no Rs, no ra => neither causal nor generalisable, ONLY Correlational [BAD Observations]

Scatter plot => Explanatory variable on X axis, Response on Y. Association might be identified, causality cannot.

Identifying relationship

- Direction: positive, negative
- Shape: Linear, curve
- Strength: Strong, weak
- Outliers

Skewness is determined by the longer tail! Modality: Peak in the distribution, unimodal, bimodal, uniform, multimodal

Boxplot: Median and IQR . Using boxplot, you can sketch back the histogram

Mean < Median : Left skewed distribution Mean > Median: Right skewed distribution Mean = Median : Symmetric

Measures of spread:

- Range = max-min
- Variance = average squared deviation from mean

$$s^2 = \frac{\sigma(x + i - x)^2}{n - 1} \quad (1)$$

- Std deviation