# Observations vs Experiments

measures specific trait but does not modify subjects

Apply a treatment and then measure the effect on the subjects

Objective of sampling is to generalize properties of the population.

Random — each member of the population have an equal probability of being picked or selected in the sample.         chance

## a.) Simple Random Sampling.

Each group of size "n" has an equal chance of being selected.

4 common ways to achieve this —

eg:- what is peoples favourite movie right now? - you ask people on your phone list. This is not going to be random.

① Convenience Sampling.
- use the results that are easy to get. (Not random)

② Systematic Sampling.
- put a population in order and select every "kth" member.

③ Stratified Sampling.          strata = layers   (stratosphere)
- In systematic sampling, there is a chance that you will miss out on a certain group.
In stratified, you divide the population into subgroups based on a characteristic, then sample each subgroup.

④ Cluster Sampling.
- divide population into groups that need not have any characteristics that are similar.
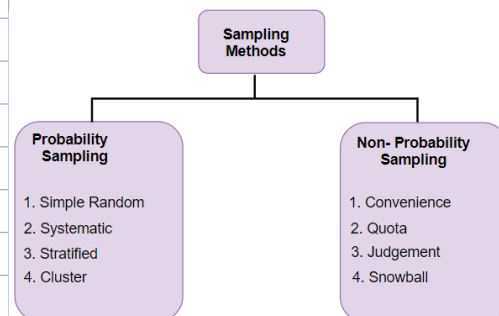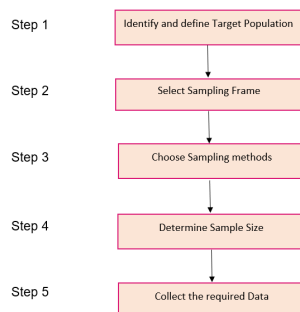- Randomly selected a certain number of clusters and then sample or collect data from the entire cluster.
   eg- Divide a class into 4 groups. and pick 2 groups randomly and then sample data.

Sampling Error —
↳ Difference in characteristics between your sample and the population.

| | |
|---|---|
| Step 1 | Identify and define Target Population |
| Step 2 | Select Sampling Frame |
| Step 3 | Choose Sampling methods |
| Step 4 | Determine Sample Size |
| Step 5 | Collect the required Data |

**Sampling Methods**

**Probability Sampling**
1. Simple Random
2. Systematic
3. Stratified
4. Cluster

**Non- Probability Sampling**
1. Convenience
2. Quota
3. Judgement
4. Snowball

# sampling and resampling!

## Sampling with replacement and without replacement.

With replacement, → same observation can repeat. If sampling 2 events
the occurence of first event does not impact other event.
covariance is zero.
   i.e samples are independent.
   eg. pick 2 samples from a set with below numbers.
   (1,1) can appear --→ repeated values $\{1, 2, 3, 4, 5, 6\}$
   $P = \frac{1}{6} \times \frac{1}{6}$.

without replacement, samples cannot repeat. covariance is non-zero. events
depends on one another.
   eg.   $P(1, 2)$       basically, when you pick the first observation, you remove
   $P = \frac{1}{6} \times \frac{1}{5}$       it from the original population.

Random - number generators are used to create random samples.
There is no such thing as "truly" random. Hence called pseudo-
random number generators (PRNG)

"seed" variable is used to initialize this PRNG algorithms.

Python random() module implements these algorithms.

Linear congruential method implemented in random() function.
$$A_{n+1} = (Z \times A_n + I) \bmod M$$

   $Z, I, M$ all constant numbers.

This process is fairly deterministic / predictable.