# Sampling Distributions (EPIB 607)

Sahir Bhatnagar and James Hanley

Department of Epidemiology, Biostatistics, and Occupational Health
McGill University

sahir.bhatnagar@mcgill.ca
https://sahirbhatnagar.com/EPIB607/

September 17, 2018

Parameters, Samples, and Statistics

# Parameters, Samples, and Statistics

- **Paramter**: An unknown numerical constant pertaining to a population/universe, or in a statistical model.
  - $\mu$: population mean      $\pi$: population proportion
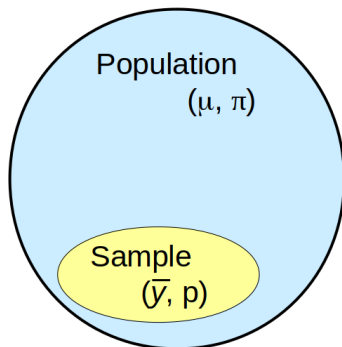
# Parameters, Samples, and Statistics

- **Paramter**: An unknown numerical constant pertaining to a population/universe, or in a statistical model.
  - ▶ $\mu$: population mean      $\pi$: population proportion
- **Statistic**: A numerical quantity calculated from a sample. The empirical counterpart of the parameter, used to *estimate* it.

# Parameters, Samples, and Statistics

- **Paramter**: An unknown numerical constant pertaining to a population/universe, or in a statistical model.
  - $\mu$: population mean        $\pi$: population proportion
- **Statistic**: A numerical quantity calculated from a sample. The empirical counterpart of the parameter, used to *estimate* it.
  - $\bar{y}$: sample mean        $p$: sample proportion

# Parameters, Samples, and Statistics

- **Paramter**: An unknown numerical constant pertaining to a population/universe, or in a statistical model.
  - ▸ $\mu$: population mean $\qquad$ $\pi$: population proportion
- **Statistic**: A numerical quantity calculated from a sample. The empirical counterpart of the parameter, used to *estimate* it.
  - ▸ $\bar{y}$: sample mean $\qquad$ $p$: sample proportion



Population
$(\mu, \pi)$

Sample
$(\bar{y}, p)$

# Examples

Proportions:

- Proportion of Earth's surface covered by water

- Proportion who saw a medical doctor last year

- Proportion of Québécois who don't have a family doctor

# Examples

**Proportions**:

- Proportion of Earth's surface covered by water

- Proportion who saw a medical doctor last year

- Proportion of Québécois who don't have a family doctor

**Means**:

- Mean depth in $n$ randomly selected ocean locations

- Mean household size in $n$ randomly selected households.

- Median number of persons under-5 in a sample of $n$ households

# Samples must be random

- The validity of inference will depend on the way that the sample was collected. If a sample was collected badly, no amount of statistical sophistication can rescue the study.

# Samples must be random

- The validity of inference will depend on the way that the sample was collected. If a sample was collected badly, no amount of statistical sophistication can rescue the study.

- Samples should be **random**. That is, there should be no systematic set of characteristics that is related to the scientific question of interest that causes some people to be more likely to be sampled than others. The simplest type of randomization selects members from the population with equal probability (a uniform distribution).

# Samples must be random

- The validity of inference will depend on the way that the sample was collected. If a sample was collected badly, no amount of statistical sophistication can rescue the study.

- Samples should be **random**. That is, there should be no systematic set of characteristics that is related to the scientific question of interest that causes some people to be more likely to be sampled than others. The simplest type of randomization selects members from the population with equal probability (a uniform distribution).

- When conducting a study, it is always better to seek statistical advice sooner rather than later. Get a statistician involved at the *planning* stage of the study... by the analysis stage, it may be too late!

# Samples must be random - No cheating!

Do not cheat by

# Samples must be random - No cheating!

**Do not cheat by**

- Taking 5 people from the *same* household to estimate
    - proportion of Québécois who don't have a family doctor
    - who saw a medical doctor last year
    - average rent

# Samples must be random - No cheating!

**Do not cheat by**

- Taking 5 people from the *same* household to estimate
  - ▶ proportion of Québécois who don't have a family doctor
  - ▶ who saw a medical doctor last year
  - ▶ average rent

- Sampling the depth of the ocean *only around Montreal* to estimate
  - ▶ proportion of Earth's surface covered by water

# Collecting data takes effort

**In general**

- The larger the sample $\rightarrow$ the more accurate the estimate (if sampling is done correctly)

# Collecting data takes effort

**In general**

- The larger the sample $\rightarrow$ the more accurate the estimate (if sampling is done correctly)

**CAVEAT**

- Collecting more data takes effort and money!
- We will also soon discover the curse of the $\sqrt{n}$

# Collecting data takes effort

**In general**

- The larger the sample $\rightarrow$ the more accurate the estimate (if sampling is done correctly)

# Collecting data takes effort

**In general**

- The larger the sample $\rightarrow$ the more accurate the estimate
  (if sampling is done correctly)

**CAVEAT**

- Collecting more data takes effort and money!
- We will also soon discover the curse of the $\sqrt{n}$

# Sampling Distributions

# Sampling Distributions

- Given a sample of $n$ observations from a population, we will be calculating estimates of the population mean, proportion, standard deviation, and various other population characteristics (parameters)

# Sampling Distributions

- Given a sample of *n* observations from a population, we will be calculating estimates of the population mean, proportion, standard deviation, and various other population characteristics (parameters)

- Prior to obtaining data, there is uncertainty as to which of all possible samples will occur

# Sampling Distributions

- Given a sample of *n* observations from a population, we will be calculating estimates of the population mean, proportion, standard deviation, and various other population characteristics (parameters)

- Prior to obtaining data, there is uncertainty as to which of all possible samples will occur

- Because of this, estimates such as $\bar{y}$ (the sample mean) will vary from one sample to another

# Sampling Distributions

- The behavior of such estimates in many samples of equal size is described by what are called **sampling distributions**

# Sampling Distributions

- The behavior of such estimates in many samples of equal size is described by what are called **sampling distributions**

- B&M definition: The sampling distribution of a statistic is the distribution of values taken by the statistic in all possible samples of the same size from the same population.

# Why are sampling distributions important?

■ They tell us how far from the target (true value of the parameter) our statistical *shot* at it (i.e. the statistic calculated form a sample) is likely to be, or, to have been.

# Why are sampling distributions important?

- They tell us how far from the target (true value of the parameter) our statistical *shot* at it (i.e. the statistic calculated form a sample) is likely to be, or, to have been.

- Thus, they are used in confidence intervals for parameters. Specific sampling distributions (based on a null value for the parameter) are also used in statistical tests of hypotheses.

# Exercise 1: How Deep is the Ocean?

- We will get a sense of what a sampling distribution is in Exercise 1

# Exercise 1: How Deep is the Ocean?

- We will get a sense of what a sampling distribution is in Exercise 1

- **CAVEAT**: This is a luxury using a toy example. In actual studies, we only get one shot!