

# lab2

April 16, 2021

```
[1]: import numpy as np
import pandas as pd
import altair as alt
```

## 1 Lab 1: Sampling design and statistical bias

In the following scenarios you'll explore through simulation how nonrandom sampling can produce datasets with statistical properties that are distorted relative to the population that the sample was drawn from. This kind of distortion is known as **bias**.

In common usage, the word 'bias' means disproportion or unfairness. In statistics, the concept has the same connotation – biased sampling favors certain observational units over others, and biased estimates are estimates that favor larger or smaller values than the truth.

This lab has you explore sampling bias. The goal is to refine your understanding about what (statistical) bias is and is not, and develop your intuition about potential mechanisms by which bias is introduced and the effect that this can have on sample statistics.

### 1.0.1 Objectives

- Simulate biased and unbiased sampling designs
  - Examine the impact of sampling bias on the sample mean
  - Apply a simple bias correction by inverse probability weighting
- 

### 1.1 Background

**Sampling design** The **sampling design** of a study refers to *the way observational units are selected* from the sampling frame (the collection of all observable units). Any design can be expressed by the probability that each unit is included in the sample. In a random sample, all units are equally likely to be included.

For example, you might want to learn about U.S. residents (population), but only be able for ethical reasons to study adults (sampling frame), and decide to do a mail survey of 2000 randomly selected addresses in each state (sampling design). (This is not a random sample of all individuals because individuals share addresses and the population sizes are different from state to state.)

**Bias** Formally, **bias** describes *the ‘typical’ deviation of a sample statistic (observed) from its population counterpart (unobserved)*.

For example, if a particular sampling design tends to produce an average measurement around 1.5 units, but the true average in the population is 2 units, then the estimate has a bias of -0.5 units. The language ‘typical’ and ‘tends to’ is important here. Estimates are rarely (almost never) perfect, so just because an estimate is off by -0.5 units for one sample doesn’t make it biased – it is only biased if it is *consistently* off under repeated sampling.

Although bias is technically a property of a sample statistic (like the sample average), it’s common to talk about a biased *sample* – this term refers to a dataset collected using a sampling design that produces biased statistics.

This is exactly what you’ll explore in this lab – the relationship between sampling design and bias.

**Simulated data** You will be simulating data in this lab. **Simulation** is a great means of exploration for the present topic *because you can control the population properties*.

When working with real data, you just have one dataset, and you don’t know any of the properties of the population or what might have happened if a different sample were collected. That makes it difficult to understand sampling variation and impossible to directly compare the sample properties to the population properties!

With simulated data, by contrast, you control how data are generated with exact precision – so by extension, you know everything there is to know about the population. In addition, repeated simulation of data makes it possible to explore the typical behavior of a particular sampling design, so you can learn ‘what usually happens’ for a particular sampling design by direct observation.

---

## 1.2 Scenario 1: unbiased samples

In this scenario you’ll compare the sample mean and the distribution of sample values for a single variable with the population mean and distribution for an unbiased sampling design.

### 1.2.1 Hypothetical population

To provide a little context to this scenario, imagine that you’re measuring eucalyptus seeds to determine their typical diameter. The cell below simulates diameter measurements for a hypothetical population of 5000 seeds; imagine that this is the total number of seeds in a small grove at some point in time.

```
[2]: # simulate seed diameters
np.random.seed(40221) # for reproducibility
population = pd.DataFrame(
    data = {'diameter': np.random.gamma(shape = 2, scale = 1/2, size = 5000),
            'seed': np.arange(5000)}
).set_index('seed')
```

```
# check first few rows
population.head(3)
```

```
[2]:      diameter
      seed
0      0.831973
1      1.512187
2      0.977392
```

**Question 1a** Calculate the mean diameter for the hypothetical population.

```
[3]: # solution

population["diameter"].mean()
```

```
[3]: 1.0189291497049837
```

**Question 1b** Calculate the standard deviation of diameters for the hypothetical population.

```
[4]: # solution

population["diameter"].std()
```

```
[4]: 0.7239297185874436
```

The cell below produces a histogram of the population values – the distribution of diameter measurements among the hypothetical population – with a vertical line indicating the population mean.

```
[5]: # base layer
base_pop = alt.Chart(population).properties(width = 400, height = 300)

# histogram of diameter measurements
hist_pop = base_pop.mark_bar(opacity = 0.8).encode(
    x = alt.X('diameter',
              bin = alt.Bin(maxbins = 20),
              title = 'Diameter (mm)',
              scale = alt.Scale(domain = (0, 6))),
    y = alt.Y('count()', title = 'Number of seeds in population')
)

# vertical line for population mean
mean_pop = base_pop.mark_rule(color='blue').encode(
    x = 'mean(diameter)'
)

# display
```

```
hist_pop + mean_pop
```

```
[5]: alt.LayerChart(...)
```

### 1.2.2 Hypothetical sampling design

Imagine that your sampling design involves collecting bunches of plant material from several locations in the grove and sifting out the seeds with a fine sieve until you obtaining 250 seeds. We'll suppose that using your collection method, any of the 5000 seeds is equally likely to be obtained, so that your 250 seeds comprise a *random sample* of the population.

We can simulate samples obtained using your hypothetical design by drawing values without replacement from the population.

```
[6]: # draw a random sample of seeds
np.random.seed(40221) # for reproducibility
sample = population.sample(n = 250, replace = False)
```

**Question 1c** Calculate the mean diameter of seeds in the simulated sample. Is it close to the population mean?

```
[7]: # solution

sample["diameter"].mean()
```

```
[7]: 0.9777218824084053
```

**Answer** *The mean diameter is close to the population mean's diameter.*

The cell below produces a histogram of the sample values, and displays it alongside the histogram of population values.

```
[8]: # base layer
base_samp = alt.Chart(sample).properties(width = 400, height = 300)

# histogram of diameter measurements
hist_samp = base_samp.mark_bar(opacity = 0.8).encode(
    x = alt.X('diameter',
              bin = alt.Bin(maxbins = 20),
              scale = alt.Scale(domain = (0, 6)),
              title = 'Diameter (mm)'),
    y = alt.Y('count()', title = 'Number of seeds in sample')
)

# vertical line for population mean
mean_samp = base_samp.mark_rule(color='blue').encode(
```

```

    x = 'mean(diameter)'
)

# display
hist_samp + mean_samp | hist_pop + mean_pop

```

```
[8]: alt.HConcatChart(...)
```

Notice that while there are some small differences, the overall shape is similar and the sample mean is almost exactly the same as the population mean. So with this sampling design, you obtained a dataset with few distortions of the population properties, and the sample mean is a good estimate of the population mean.

### 1.2.3 Assessing bias

You may wonder: *does that happen all the time, or was this just a lucky draw?* This question can be answered by simulating a large number of samples to see whether the undistorted representation of the population is typical for this sampling design. To simplify life a little, let's focus on whether the sample mean is usually accurate.

The cell below estimates the bias of the sample mean by:

- drawing 1000 samples of size 300;
- storing the sample mean from each sample;
- computing the average difference between the sample means and the population mean.

```

[9]: np.random.seed(40221) # for reproducibility

# number of samples to simulate
nsim = 1000

# storage for the sample means
samp_means = np.zeros(nsim)

# repeatedly sample and store the sample mean
for i in range(0, nsim):
    samp_means[i] = population.sample(n = 250, replace = False).mean()

```

The bias of the sample mean is its average distance from the population mean. We can estimate this using our simulation results as follows:

```

[10]: # bias
      samp_means.mean() - population.diameter.mean()

```

```
[10]: -0.0012458197406362004
```

So the average error observed in 1000 simulations was about 0.001 mm! This suggests that the sample mean is *unbiased*: on average, there is no error. Therefore, at least with respect to estimating

the population mean, random samples appear to be *unbiased samples*.

However, **unbiasedness does not mean that you won't observe estimation error**. There is a natural amount of variability from sample to sample, because in each sample a different collection of seeds is measured.

The cell below plots a histogram representing the distribution of values of the sample mean across the 1000 samples you simulated (this is known as the *sampling distribution* of the sample mean). It shows a peak right at the population mean (blue vertical line) but some symmetric variation to either side – most values are between about 0.93 and 1.12.

```
[11]: # plot the simulated sampling distribution
sampling_dist = alt.Chart(pd.DataFrame({'sample mean': samp_means})).mark_bar().
    ↪ encode(
        x = alt.X('sample mean', bin = alt.Bin(maxbins = 30), title = 'Value of_
        ↪ sample mean'),
        y = alt.Y('count()', title = 'Number of simulations')
    )

sampling_dist + mean_pop
```

```
[11]: alt.LayerChart(...)
```

---

## 1.3 Scenario 2: biased sampling

In this scenario, you'll use the same hypothetical population of eucalyptus seed diameter measurements and explore the impact of a biased sampling design.

### 1.3.1 Hypothetical sampling design

In the first design, you were asked to imagine that you collected and sifted plant material to obtain seeds. Suppose you didn't know that the typical seed is about 1mm in diameter and decided to use a sieve that is a little too coarse, tending only to sift out larger seeds and letting smaller seeds pass through. As a result, small seeds have a lower probability of being included in the sample and large seeds have a higher probability of being included in the sample.

This kind of sampling design can be described by assigning differential *sampling weights*  $w_1, \dots, w_N$  to each observation. The cell below defines a `weight_fn` that calculates a weight  $w_i$  between 0 and 1 according to diameter, so that larger diameters have larger weights and are more likely to be sampled.

```
[12]: # inclusion weight as a function of seed diameter
def weight_fn(x, r = 2, c = 2):
    out = 1/(1 + np.e**(-r*(x - c)))
    return out
```

```

# create a grid of values to use in plotting the function
grid = np.linspace(0, 6, 100)
weight_df = pd.DataFrame(
    {'seed diameter': grid,
     'weight': weight_fn(grid)}
)

# plot of inclusion probability against diameter
weight_plot = alt.Chart(weight_df).mark_area(opacity = 0.3, line = True).encode(
    x = 'seed diameter',
    y = 'weight'
).properties(height = 100)

# show plot
weight_plot

```

[12]: alt.Chart(...)

The actual probability that a seed is included in the sample – its **inclusion probability** – is proportional to the sampling weight. These inclusion probabilities  $\pi_i$  can be calculated by normalizing the weights  $w_i$ :

$$\pi_i = \frac{w_i}{\sum_i w_i}$$

It may help you to picture how the weights will be used in sampling to line up this plot with the population distribution. In effect, we will sample more from the right tail of the population distribution, where the weight is nearest to 1.

[13]: hist\_pop & weight\_plot

[13]: alt.VConcatChart(...)

The following cell draws a sample with replacement from the hypothetical seed population *with seeds weighted according to the inclusion probability given by the function above*.

```

[14]: # assign inclusion probability to each seed
population['inclusion_prob'] = weight_fn(population.diameter)/
    ↪ (weight_fn(population.diameter)).sum()

# draw weighted sample
np.random.seed(40721)
sample = population.sample(n = 250, replace = False, weights =
    ↪ 'inclusion_prob').drop(columns = 'inclusion_prob')

```

**Question 2a** Calculate the mean diameter of seeds in the simulated sample. Is it close to the population mean?

```
[15]: # solution
      sample["diameter"].mean()
```

```
[15]: 1.8117214583859176
```

**Answer** *The mean diameter is not close to the population mean's diameter.*

**Question 2b** Show side-by-side plots of the distribution of sample values and the distribution of population values, with vertical lines indicating the corresponding mean on each plot. (*Hint: copy the cell that produced this plot in scenario 1.*) Does the distribution of diameters of seeds in the sample seem to accurately reflect the population?

```
[16]: # solution
      sample["diameter"].mean()

      # base layer

      base_samp = alt.Chart(sample).properties(width = 400, height = 300)

      # histogram of diameter measurements

      hist_samp = base_samp.mark_bar(opacity = 0.8).encode(
        x = alt.X("diameter",
                  bin = alt.Bin(maxbins = 20),
                  scale = alt.Scale(domain = (0, 6)),
                  title = 'Diameter (mm)'),
        y = alt.Y('count()', title = 'Sample Number of Seeds')
      )

      # vertical line for population mean

      mean_samp = base_samp.mark_rule(color='blue').encode(x = 'mean(diameter)')

      # display

      hist_samp + mean_samp | hist_pop + mean_pop
```

```
[16]: alt.HConcatChart(...)
```

**Answer** *The depiction of the diameter of seeds does not accurately reflect the population*



### 1.3.2 Assessing bias

Here you'll mimic the simulation done in scenario 1 to assess the bias of the sample mean under this new sampling design.

**Question 2c** Investigate the bias of the sample mean by:

- drawing 1000 samples with observations weighted by inclusion probability;
- storing the sample mean from each sample;
- computing the average difference between the sample means and the population mean.

(*Hint*: copy the cell that performs this simulation in scenario 1, and be sure to adjust the sampling step to include `weights = ...` with the appropriate argument.)

```
[17]: # solution
np.random.seed(40221) # for reproducibility

# number of samples to simulate

nsim = 1000

# storage for the sample means

samp_means = np.zeros(nsim)

# repeatedly sample and store the sample mean

for i in range(0, nsim):
    samp_means[i] = population.sample(n = 250,
                                     replace = False,
                                     weights = "inclusion_prob").drop(columns = "inclusion_prob").mean()

# bias

samp_means.mean() - population.diameter.mean()
```

```
[17]: 0.7722765557571376
```

**Question 2d** Does this sampling design seem to introduce bias? If so, does the sample mean tend to over-estimate or under-estimate the population mean?

**Answer** *There is bias introduced, as the population mean is over-estimated*

---

## 1.4 Scenario 3

In this scenario, you'll explore sampling from a population with group structure – frequently bias can arise from inadvertent uneven sampling of groups within a population.

### 1.4.1 Hypothetical population

Suppose you're interested in determining the average beak-to-tail length of red-tailed hawks to help differentiate them from other hawks by sight at a distance. Females and males differ slightly in length – females are generally larger than males. The cell below generates length measurements for a hypothetical population of 3000 females and 2000 males.

```
[18]: # for reproducibility
np.random.seed(40721)

# simulate hypothetical population
population = pd.DataFrame(
    data = {'length':np.random.normal(loc =57.5,scale = 3,size = 3000),
           'sex':np.repeat('female', 3000)}
).append(
    pd.DataFrame(
        data = {'length':np.random.normal(loc = 50.5,scale = 3,size = 2000),
               'sex':np.repeat('male', 2000)}
    )
)

# preview
population.groupby('sex').head(2)
```

```
[18]:      length      sex
0  53.975230  female
1  60.516768  female
0  53.076663   male
1  49.933166   male
```

The cell below produces a histogram of the lengths in the population overall (bottom panel) and when distinguished by sex (top panel).

```
[19]: base = alt.Chart(population).properties(height = 200)

hist = base.mark_bar(opacity = 0.5, color = 'red').encode(
    x = alt.X('length',
              bin = alt.Bin(maxbins = 40),
              scale = alt.Scale(domain = (40, 70)),
              title = 'length (cm)'),
    y = alt.Y('count()',
              stack = None,
```

```

        title = 'number of birds')
    )

hist_bysex = hist.encode(color = 'sex').properties(height = 100)

hist_bysex & hist

```

```
[19]: alt.VConcatChart(...)
```

The population mean – average length of both female and male red-tailed hawks – is shown below.

```
[20]: # population mean
      population.mean()
```

```
[20]: length    54.737717
      dtype: float64
```

First try drawing a random sample from the population:

```
[21]: # for reproducibility
      np.random.seed(40821)

      # randomly sample
      sample = population.sample(n = 300, replace = False)
```

**Question 3a** Do you expect that the sample will contain equal numbers of male and female hawks? Think about this for a moment (you don't have to provide a written answer), and then compute the proportions of individuals in the sample of each sex.

(Hint: group by sex, use `.count()`, and divide by the sample size. Be sure to rename the output column appropriately, as the default behavior produces a column called `length`.)

```
[22]: # solution

SexSmp = sample.groupby('sex').count().rename(columns={'length': 'Gender %'})/300
SexSmp
```

```
[22]:      Gender %
sex
female  0.596667
male    0.403333
```

The sample mean is shown below, and is fairly close to the population mean. This should be expected, since you already saw in scenario 1 that random sampling is an unbiased sampling design with respect to the mean.

```
[23]: # solution
      sample["length"].mean()
```

[23]: 54.95210318784532

### 1.4.2 Biased sampling

Let's now consider a biased sampling design. Usually, length measurements are collected from dead specimens collected opportunistically. Imagine that male mortality is higher, so there are better chances of finding dead males than dead females. Suppose in particular that specimens are five times as likely to be male; to represent this situation, we'll assign sampling weights of 5/6 to all male hawks and weights of 1/6 to all female hawks.

```
[24]: def weight_fn(sex, p = 5/6):
        if sex == 'male':
            out = p
        else:
            out = 1 - p
        return out

weight_df = pd.DataFrame(
    {'length': [50.5, 57.5],
     'weight': [5/6, 1/6],
     'sex': ['male', 'female']})

wt = alt.Chart(weight_df).mark_bar(opacity = 0.5).encode(
    x = alt.X('length', scale = alt.Scale(domain = (40, 70))),
    y = alt.Y('weight', scale = alt.Scale(domain = (0, 1))),
    color = 'sex'
).properties(height = 70)

hist_bysex & wt
```

[24]: alt.VConcatChart(...)

**Question 3b** Draw a weighted sample from the population using the weights defined by `weight_fn`, and compute the sample mean.

```
[25]: # for reproducibility
np.random.seed(40821)

# assign weights
population["weight"] = population.sex.aggregate(func=weight_fn)

# randomly sample
sample = population.sample(n=300, replace=False, weights="weight").
    ↳ drop(columns="weight")
```

```
# compute mean
sample.mean()
```

```
[25]: length      51.887106
      dtype: float64
```

**Question 3c** Investigate the bias of the sample mean by:

- drawing 1000 samples with observations weighted by `weight_fn`;
- storing the sample mean from each sample;
- computing the average difference between the sample means and the population mean.

```
[26]: # solution
      np.random.seed(40221) # for reproducibility

      # number of samples to simulate
      nsim = 1000

      # storage for the sample means
      samp_means = np.zeros(nsim)

      # repeatedly sample and store the sample mean
      for i in range(0, nsim):
          samp_means[i] = population.sample(n = 300, replace = False, weights = ↵
          ↵ "weight").drop(columns = "weight").mean()

      # bias
      samp_means.mean() - population.length.mean()
```

```
[26]: -2.5720649894415715
```

**Question 3d** Reflect a moment on your simulation result in question 3c. If *female* mortality is higher and specimens for measurement are collected opportunistically, as described in this sampling design, do you expect that the average length in the sample will be an underestimate or an overestimate of the population mean? Explain why in 1-2 sentences.

**Answer** \*The average length in the sample will be an overestimate of the population mean. This is shown through the graph indicating female hawks have a larger length, on average, in beak to tail measurements than males, resulting in an overestimate because a larger number of female hawks are being sampled.

---

## 1.5 Bias correction

*What can be done if a sampling design is biased? Is there any remedy?*

You've seen some examples above of how bias can arise from a sampling mechanism in which units have unequal chances of being selected in the sample. Ideally, we'd work with random samples all the time, but that's not very realistic in practice. Fortunately, biased sampling is not a hopeless case – **it is possible to apply bias corrections if you have good information about which individuals were more likely to be sampled.**

To illustrate how this would work, let's revisit scenario 2 – sampling larger eucalyptus seeds more often than smaller ones. Imagine you realize the mistake and conduct a quick experiment with your sieve to determine the proportion of seeds of each size that pass through, and use this to estimate the inclusion probabilities. (To simplify this exercise, we'll just use sampling weights we defined to calculate the actual inclusion probabilities.)

The cell below generates the population and sample from scenario 2 again:

```
[27]: # simulate seed diameters
np.random.seed(40221) # for reproducibility
population = pd.DataFrame(
    data = {'diameter': np.random.gamma(shape = 2, scale = 1/2, size = 5000),
            'seed': np.arange(5000)}
).set_index('seed')

# probability of inclusion as a function of seed diameter
def weight_fn(x, r = 2, c = 2):
    out = 1/(1 + np.e**(-r*(x - c)))
    return out

# assign inclusion probability to each seed
population['samp_weight'] = weight_fn(population.diameter)

# draw weighted sample
np.random.seed(40721)
sample = population.sample(n = 250, replace = False, weights = 'samp_weight')
```

The sample mean and population mean you calculated earlier are shown below:

```
[28]: # print sample and population means
pd.Series({'sample mean': sample.diameter.mean(), 'population mean': population.
    ↪diameter.mean()})
```

```
[28]: sample mean      1.811721
      population mean  1.018929
      dtype: float64
```

We can obtain an unbiased estimate of the population mean by computing a *weighted average* of the diameter measurements instead of the sample average after weighting the measurements in inverse proportion to the sampling weights:

$$\text{weighted average} = \sum_{i=1}^{250} \underbrace{\left( \frac{w_i^{-1}}{\sum_{j=1}^{250} w_j^{-1}} \right)}_{\text{bias adjustment}} \times \text{diameter}_i$$

This might look a little complicated, but the idea is simple – the weighted average corrects for bias by simply up-weighting observations with a lower sampling weight and down-weighting observations with a higher sampling weight.

The cell below performs this calculation.

```
[29]: # compute bias adjustment
sample['bias_adjustment'] = (sample.samp_weight**(-1))/(sample.
    ↪ samp_weight**(-1)).sum()

# weight diameter measurements
sample['weighted_diameter'] = sample.diameter*sample.bias_adjustment

# sum to compute weighted average
sample.weighted_diameter.sum()
```

```
[29]: 1.0139201948221341
```

Notice that the weighted average successfully corrected for the bias:

```
[30]: # print sample and population means
pd.Series({'sample mean': sample.diameter.mean(),
    'weighted average': sample.weighted_diameter.sum(),
    'population mean': population.diameter.mean()})
```

```
[30]: sample mean      1.811721
      weighted average  1.013920
      population mean  1.018929
      dtype: float64
```

## 1.6 Takeaways

These simulations illustrate through a few simple examples that random sampling – a sampling design where each unit is equally likely to be selected – produces unbiased sample means. That means that ‘typical samples’ will yield sample averages that are close to the population value. By contrast, deviations from random sampling tend to yield biased sample averages – in other words, nonrandom sampling tends to distort the statistical properties of the population in ways that can produce misleading conclusions (if uncorrected).

Here are a few key points to reflect on:

- bias is not a property of an individual sample, but of a *sampling design*

- unbiased sampling designs tend to produce faithful representations of populations
    - but there are no guarantees for individual samples
  - if you hadn't known the population distributions, there would have been no computational method to detect bias
    - in practice, it's necessary to *reason* about whether the sampling design is sound
  - the sample statistic (sample mean) was only reliable when the sampling design was sound
    - the quality of data collection is arguably more important for reaching reliable conclusions than the choice of statistic or method of analysis
- 

## 1.7 Submission Checklist

1. Save file to confirm all changes are on disk
2. Run *Kernel > Restart & Run All* to execute all code from top to bottom
3. Save file again to write any new output to disk
4. Select *File > Download* (should save as .ipynb)
5. Submit to Gradescope