



# Homework 6

PSTAT 5A: Spring 2023, with Ethan P. Marzban

## **i** Instructions

- Please submit your work to Gradescope by no later than **11:59pm on Tuesday, May 16**. As a reminder, late homework will not be accepted.
- Recall that you will be asked to upload a **single** PDF containing your work for *both* the programming and non-programming questions to Gradescope.
  - You can merge PDF files using either Adobe Acrobat, or using adobe's online PDF merger at [this link](#).

## Problem 1: Airplanes in the Night Sky

An aviation enthusiast is interested in determining the proportion of all US citizens that have never flown on a plane. She takes a representative sample of 210 US Citizens, and finds out that 18% of these people have never flown on a plane.

- a. Is the value of 18% a population parameter or an observed instance of a sample statistic?

**Solution:** Because the 18% is associated with a particular *sample* of people, it is an observed instance of a sample statistic.

- b. Construct a 95% confidence interval for the true proportion of US citizens that have never flown on a plane before, and interpret your interval in the context of the problem.

**Solution:** A 95% confidence interval for a population proportion  $p$ , based on a sample of size  $n$  that resulted in an observed sample proportion  $\hat{p}$ , will take the form

$$\hat{p} \pm 1.96 \cdot \sqrt{\frac{\hat{p} \cdot (1 - \hat{p})}{n}}$$

In this case we have  $\hat{p} = 0.18$  and  $n = 210$ , meaning our confidence interval takes the form

$$(0.18) \pm 1.96 \cdot \sqrt{\frac{(0.18) \cdot (1 - 0.18)}{210}} \approx (0.18) \pm (0.052) = [0.128, 0.232]$$

The interpretation of this interval is as follows: we are 95% confident that the true proportion of US citizens that have never flown on a plane before lies between 12.8% and 23.2%.

- c. If we were to construct an 82% confidence interval for the true proportion of US citizens that have never flown on a plane before, would we expect this interval to be wider or narrower

than the interval constructed in part (c)? **Do not construct the interval first.**

**Solution:** Recall from lecture that higher confidence levels correspond to wider intervals; conversely, lower confidence intervals correspond to narrower intervals. As such, we would expect an 82% confidence interval to be **narrower** than a 95% confidence interval.

- d. Now, construct an 82% confidence interval for the true proportion of US citizens that have never flown on a plane before. Comment on whether this agrees with your answer to part (b) or not.

**Solution:** Our first task is to find the value of  $z_\alpha$  that corresponds to an 82% confidence level. Following the formula provided in class, we need to look up the

$$\left(\frac{1 - 0.82}{2}\right) \times 100$$

i.e. the 9th percentile of the standard normal distribution. From a standard normal table we see that  $P(Z \leq -1.34) = 0.0901$  which is close enough to 0.09 for us to take  $z_\alpha = 1.34$ . As such, our confidence interval then becomes

$$(0.18) \pm 1.34 \cdot \sqrt{\frac{(0.18) \cdot (1 - 0.18)}{210}} \approx (0.18) \pm (0.036) = [0.144, 0.216]$$

Indeed, this interval is narrower than the interval we found in part (b) above.

## Problem 2: Scores

In a particular iteration of PSTAT 5A, scores on the final exam had an average of 89 and a standard deviation of 40. The exact distribution of scores is, however, unknown. Suppose a representative sample of 100 students is taken, and the average final exam score of these 100 students is recorded.

- a. Identify the population.

**Solution:** The population is the set of all students in the aforementioned iteration of PSTAT 5A.

- b. Identify the sample.

**Solution:** The sample is the 100 students that were selected.

- c. Define the parameter of interest. Use the notation discussed in Lecture 12.

**Solution:** We use  $\mu$  to denote population means; as such, let  $\mu$  denote the true average final exam score of PSTAT 5A students.

- d. Define the random variable of interest. Use the notation discussed in Lecture 12.

**Solution:** We use  $\bar{X}$  to denote sample means; as such, let  $\bar{X}$  denote the average final exam score of 100 randomly-selected students from PSTAT 5A.

- e. What is the sampling distribution of the random variable you defined in part (d) above? Be sure to check any conditions that might need to be checked!

**Solution:** The first question we ask ourselves is: is the population normally distributed? The answer is no. As such, we then ask ourselves: is the sample size greater than 30? The answer is yes. As such, we finally ask ourselves: is the population standard deviation known? The answer is yes. Hence,  $\bar{X}$  will be normally distributed; specifically,

$$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma}{\sqrt{n}}\right) \sim \mathcal{N}\left(89, \frac{40}{\sqrt{100}}\right) \sim \mathcal{N}(89, 4)$$

- f. What is the approximate probability that the average score of these 100 students lies within 5 points of the true average score of 89?

**Solution:** We seek  $\mathbb{P}(84 \leq \bar{X} \leq 94)$ . As such, we compute

$$\begin{aligned} \mathbb{P}(84 \leq \bar{X} \leq 94) &= \mathbb{P}(\bar{X} \leq 94) - \mathbb{P}(\bar{X} \leq 84) \\ &= \mathbb{P}\left(\frac{\bar{X} - 89}{4} \leq \frac{94 - 89}{4}\right) - \mathbb{P}\left(\frac{\bar{X} - 89}{4} \leq \frac{84 - 89}{4}\right) \\ &= \mathbb{P}(Z \leq 1.25) - \mathbb{P}(Z \leq -1.25) = 0.8944 - 0.1056 = 78.88\% \end{aligned}$$

### Problem 3: An Apple a Day

Quinn is interested in performing inference on the average weight of Granny Smith apples in the Santa Barbara location of *Bristol Farms*. To that end, he takes a representative sample of 52 apples; the mean weight of his sample was 83g and the standard deviation of weights in his sample was 17g.

- a. Identify the population.

**Solution:** The population is the set of all apples at the Santa Barbara location of \*Bristol Farms\*.

- b. Identify the sample.

**Solution:** The sample is the set of 52 apples Quinn selected.

- c. Define the parameter of interest. Use the notation discussed in Lecture 12.

**Solution:** We let  $\mu$  denote the true average weight of Granny Smith apples at the Santa Barbara location of \*Bristol Farms\*.

d. Define the random variable of interest. Use the notation discussed in Lecture 12.

**Solution:** We let  $\bar{X}$  denote the average weight of a sample of 52 Granny Smith apples, taken from the Santa Barbara location of \*Bristol Farms\*.

e. What distribution do we use to construct confidence intervals for the true average weight of a Granny Smith apple at the Santa Barbara location of *Bristol Farms*?

**Solution:**

- Is the population normally distributed? No.
- Is the sample size greater than 30? Yes.
- Is the population standard deviation known? No, only the sample standard deviation.

As such, we use the  $t$  distribution with  $n - 1 = 52 - 1 = 51$  degrees of freedom; i.e. the  $t_{51}$  distribution.

f. Construct a 95% confidence interval for the true average weight of a Granny Smith apple at the Santa Barbara location of *Bristol Farms*.

**Solution:** Our confidence interval will be of the form

$$\bar{x} \pm t_{51, \alpha} \cdot \frac{s}{\sqrt{51}}$$

Here,  $\bar{x} = 83$  and  $s = 17$ . Now, the  $t$ -table does not actually have a row for 51 degrees of freedom; as such, we can either obtain an approximate value by simply using 50 degrees of freedom (which gives us a value of  $t_{51, \alpha} \approx 2.01$ ), or we can simply use Python (which also gives us a value of around 2.01). As such, our confidence interval becomes

$$83 \pm (2.01) \cdot \frac{17}{\sqrt{52}} = [78.26147, 87.73853]$$

## Problem 4: Programming

In this problem, we will build off of the material in Lab 06 to finally reproduce the demo performed in Lecture 10.

### Part (a): Extending Lists

It is often necessary to be able *extend* a list of values. As an example, consider a variable  $x = [1, 2, 3]$ . To append a value of 4 to  $x$  (i.e. to redefine  $x$  as  $[1, 2, 3, 4]$ ) we could simply redefine  $x$  to be  $[1, 2, 3, 4]$ , however for very large or otherwise complicated lists this may not be

feasible. Thankfully, Python has a built-in function to append to a list, called `append()`. To get a sense of how it works, let's do an example.

### ! Task 1

- a. Define a variable `x` that is assigned a value of `[1, 2, 3, 4, 5]`.
- b. Run `x.append(6)` **exactly once**. Do **not** run this cell multiple times.
- c. Call `x`, and see if `x` has successfully been extended.

### 💡 Solutions

#### Part (a)

```
1 x = [1, 2, 3, 4, 5]
```

#### Part (b)

```
1 x.append(6)
```

#### Part (c)

```
1 x
```

```
[1, 2, 3, 4, 5, 6]
```

### 🔥 Caution

Using command `<list>.append()` function **mutates** the `<list>` object. That is, it physically changes the structure of `<list>` (specifically by adding an additional element at the end). What this means is that each time you run the command `<list>.append()`, a new element gets appended to `<list>`. This is why it is important you only run `append()` statements *once*.

For those of you who are curious, there is a *different* `append()` function (from the `numpy` module) that does *not* mutate the list/array being appended. But, for the purposes of this class, you do not need to know how to use it.

### Part (b): Combining Loops and List Extensions

Let's combine our knowledge of `for` loops (from Lab 6) with our knowledge of list extensions (from part (a) above).

## ! Task 2

Consider the following game: roll a fair 12-sided die. If the number showing is between 1 and 4 inclusive we win a dollar; if the number showing is between 5 and 8 inclusive we win 2 dollars; if the number showing is between 9 and 12 inclusive we lose a dollar.

- Simulate playing one round of this game. As a hint: you will need to use a conditional expression somewhere in your code.
- Now, simulate playing 100 rounds of this game, and store the amount you win on each round in a list called `amount_won`. Here is a sample template you can use:

```
1 amount_won = [] # initialize a blank list to store the amount won
2
3 for <fill this in>:
4     if <something>:
5         amount_won.append(<something>)
6     elif <something>:
7         amount_won.append(<something>)
8     else:
9         amount_won.append(<something>)
```

Display the first 10 elements of the `amount_won` list.

## 💡 Solutions

### Part (a)

```
1 import random as rnd
2 import numpy as np
3
4 roll = rnd.choices(np.arange(0, 13, 1))
5
6 if 1 <= roll[0] <= 4:
7     print("Gained a dollar")
8 elif 5 <= roll[0] <= 8:
9     print("Gained two dollars")
10 else:
11     print("Lost a dollar")
```

Lost a dollar

### Part (b)

```

1 amount_won = []
2
3 for k in np.arange(0, 100):
4     roll = rnd.choices(np.arange(0, 13, 1))
5     if 1 <= roll[0] <= 4:
6         amount_won.append("Gained a dollar")
7     elif 5 <= roll[0] <= 8:
8         amount_won.append("Gained two dollars")
9     else:
10        amount_won.append("Lost a dollar")
11
12 amount_won[0:10]

```

```

['Gained a dollar',
 'Gained two dollars',
 'Gained two dollars',
 'Lost a dollar',
 'Lost a dollar',
 'Gained two dollars',
 'Lost a dollar',
 'Gained a dollar',
 'Gained a dollar',
 'Gained two dollars']

```

### Part (c): Creating a Population

We are almost in a position to fully understand the Lab 11 demo now! There is just perhaps one bit of code that we haven't explored fully yet:

#### ! Task 3

Notice that in code cell 4 of the Lab 10 demo, the following code was run:

```

1 x = rnd.choices(['Positive', 'Negative'],
2 weights = [0.035, 1 - 0.035], k = samp_size)

```

In a Markdown cell, explain *in words* what this code is doing. Specifically, think about how this code might be used to simulate drawing from a population of cats, some of which are FIV-positive. In your description, you should also mention what the population proportion (of FIV-positive cats) is in this particular simulation.

## 💡 Solutions

What this code does is simulate drawing a cat from a population of cats where 3.5% are FIV-positive. In practice, Python is selecting the element 'Positive' with *weight* (i.e. probability) 0.035, and 'Failure' with weight  $1 - 0.035$ .

### Part (d): Your Turn!

Alright, between Lab 6 and the above tasks we should have enough information to fully understand - and reproduce! - the Lab 10 demo.

## ! Task 4

In a particular country, 12% of residents live below the poverty line. Simulate taking 10,000 samples of size  $n = 250$  each, computing the proportion of people in each sample that live below the poverty line, and storing these 10,000 sample proportions in a list. Plot a histogram of the list of observed sample proportions, and overlay the true density curve of  $\hat{P}$ , the proportion of people in a representative sample of size  $n = 250$  people that live below the poverty line.

## 💡 Solutions

First, we generate the 10,000 sample proportions:

```
1 np.random.seed(10) # you don't need to do this
2
3 props = []
4 samp_size = 250
5
6 for b in np.arange(0, 10000):
7     stored_sample = rnd.choices(['Below', 'Above'],
8                                 weights = [0.12, 1 - 0.12],
9                                 k = samp_size)
10
11     count = 0
12     for k in stored_sample:
13         if k == 'Below':
14             count += 1
15
16     props.append(count / samp_size)
```

The first 10 elements of the props vector looks like:

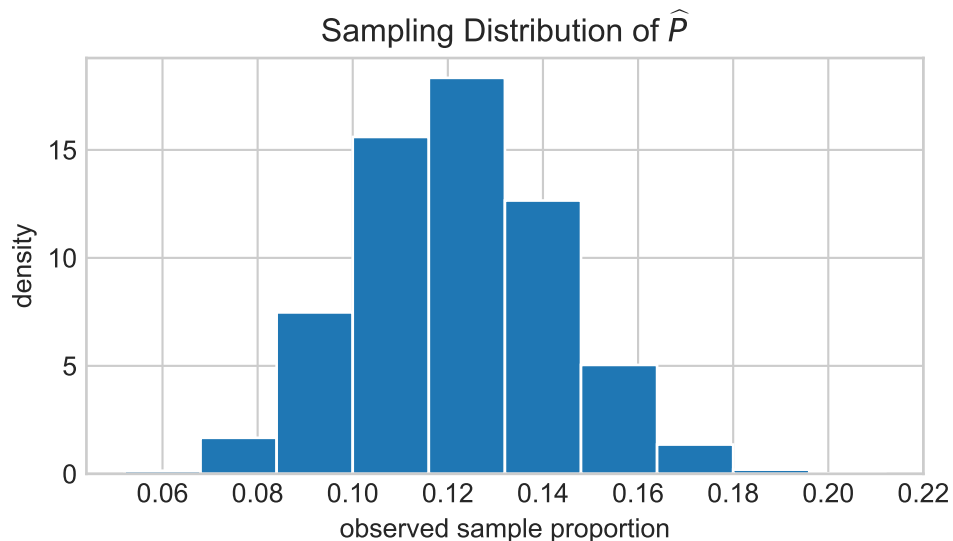


```
1 props[0:10]
```

```
[0.108, 0.132, 0.112, 0.12, 0.144, 0.144, 0.108, 0.14, 0.124, 0.132]
```

To generate a histogram of these 10,000 sample proportions, we use

```
1 %matplotlib inline
2 import matplotlib
3 import matplotlib.pyplot as plt
4 plt.style.use('seaborn-v0_8-whitegrid')
5
6 plt.figure(figsize=(5.625, 2.8125)) # optional
7 plt.hist(props, edgecolor = "white", density = True);
8 plt.xlabel("observed sample proportion");
9 plt.ylabel("density");
10 plt.title("Sampling Distribution of  $\widehat{P}$ ");
```



Finally, we need to overlay the appropriate density curve. Because the success-failure conditions are met (I leave it to you to check this), we know that

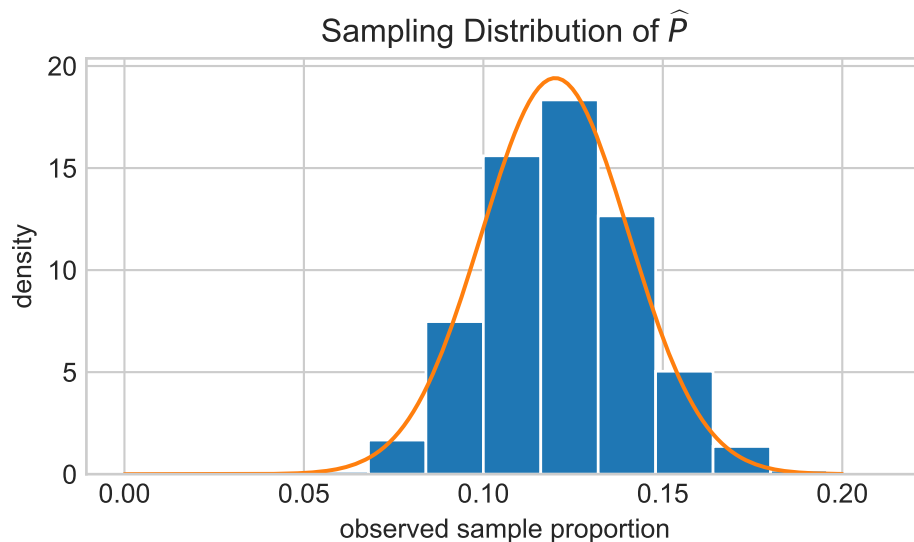
$$\widehat{P} \sim \mathcal{N}\left(p, \sqrt{\frac{p(1-p)}{250}}\right)$$

We also know the true proportion  $p$ ; that is,  $p = 0.12$ . Hence,

$$\widehat{P} \sim \mathcal{N}\left(0.12, \sqrt{\frac{0.12 \cdot (1 - 0.12)}{250}}\right)$$

the density of which can be accessed using the `scipy.stats.norm.pdf()` function.

```
1 import scipy.stats as sps
2
3 x = np.linspace(0, 0.2, 150)
4
5 plt.figure(figsize=(5.625, 2.8125)) # optional
6 plt.hist(props, edgecolor = "white", density = True);
7 plt.xlabel("observed sample proportion");
8 plt.ylabel("density");
9 plt.title("Sampling Distribution of  $\widehat{P}$ ");
10 plt.plot(x,
11          sps.norm.pdf(x, 0.12, np.sqrt(0.12 * (1 - 0.12) / samp_size))
12          );
```



### Problem 5: One More Programming Question

I would also like to quickly introduce a function that might make our lives easier when it comes to computing the value of  $z_\alpha$  in a confidence interval formula. Specifically, it is the `scipy.stats.norm.ppf()` function, which computes the percentiles of the normal distribution. For instance,

```
1 import scipy.stats as sps
2 sps.norm.ppf(0.025)
```

```
-1.9599639845400545
```

### ! Task 5

In Problem 1(d) above, you needed to find the value of  $z_\alpha$  corresponding to an 82% confidence level. Use the `scipy.stats.ppf()` function to confirm the result you obtained from the standard normal table.

### i Solutions

```
1  sps.norm.ppf((1 - 0.82)/2)
-1.340755033690216
```

We can analogously find percentiles of the  $t_k$  distribution, using `scipy.stats.t.ppf(p, k)`.

### ! Task 6

- Use Python to find the 2.5<sup>th</sup> percentile of the  $t$  distribution with 20 degrees of freedom.
- Use Python to find the 2.5<sup>th</sup> percentile of the  $t$  distribution with 10 degrees of freedom.
- Use Python to find the 2.5<sup>th</sup> percentile of the  $t$  distribution with 1000 degrees of freedom.

### 💡 Solutions

#### Part (a)

```
1  sps.t.ppf(0.025, 20)
-2.085963447265837
```

#### Part (b)

```
1  sps.t.ppf(0.025, 10)
-2.2281388519649385
```

#### Part (c)

```
1  sps.t.ppf(0.025, 1000)
-1.9623390808264078
```