

population

everything that  
you're interested  
in



sample

the part of the  
population about  
which you have  
collected data.

Not measuring the entire population,

⇒ there will be chance error in the sample.

How big will this chance error be?

What does it depend on?

Somewhat unintuitively, - depends mostly on the size  
of the sample, much less  
on the size of the population.

Start by looking at the chance error in percentage.

(population is in different categories; look at the % of  
the sample that's in a particular category).

AMS 5 students

79	fresh
83	soph.
24	juniors
16	seniors

202 total.

consider taking a simple random sample of this population.

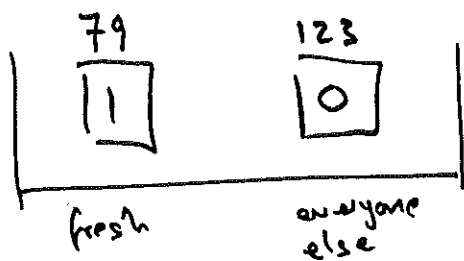
Sampling without replacement if sample size  $\ll$  population size this is ok.

sample of size 10

Expected % ~~is~~ fresh in our sample?

Size of the chance variation?

For a simple random sample, the expected value for the sample % equals the population %.



$$\text{Mean of box} = \frac{79}{202} = 0.39$$

$$\text{SD}_{\text{box}} = (1-0) \sqrt{\frac{79}{202} \times \frac{123}{202}} = 0.488.$$

Sample of size 10.

$$\text{Expected \# fresh} = 10 \times 0.39 = 3.9$$

$$\text{SE} = \sqrt{10} \times 0.488 \approx 1.5$$

expect  $3.9 \pm 1.5$  fresh.

What about %?

$$\text{Expected \% of fresh.} = \frac{\text{expected \#}}{\text{sample size}} \times 100$$

$$= \frac{3.9}{10} \times 100 = 39\%$$

( $\equiv$  % in the population)

$$\text{Chance error.} = \frac{SE}{\text{sample size}} \times 100$$

$$= \frac{1.5}{10} \times 100 = 15\%$$

We expect the sample to be  $39\% \pm 15\%$  fresh.

---

Take a sample of size 40.

$$\text{Expected \# fresh.} = \underset{\substack{\text{sample} \\ \text{size}}}{40} \times \underset{\substack{\text{mean} \\ \text{of box}}}{0.39}$$

$$\text{Expected \% fresh} = \frac{\text{Expected \#}}{\text{sample size}} \times 100.$$

$$= \frac{40 \times 0.39}{40} \times 100 = 39\%$$

chance error.

$$SE \# \text{ fresh} = \sqrt{40} \times 0.488.$$

$$\text{chance error in \% fresh} = \frac{\sqrt{40} \times 0.488}{40} \times 100$$

$$= \frac{0.488 \cdot}{\sqrt{40}} \times 100$$

$$= 7.7\%$$

Sample of size 40,  
 $39 \pm 7.7\%$  fresh.

Sample of size 10  
 $39 \pm 15\%$  fresh

Sample size has increased by  
a factor of 4.

As sample size  
increases

chance error has decreased by  
a factor of 2.

The chance error  
in \% decreases.



this is true irrespective  
of the size of the population,  
provided that the sample is small  
enough that we can consider the sample  
to be approximately drawn with replacement.

What happens when the sample is a significant fraction of the population?

The contents of the box get smaller on each draw  $\Rightarrow$  slightly less variability.

SE when drawing without replacement = SE when drawing with replacement  $\times$  correction factor

$$\text{correction factor} = \sqrt{\frac{\# \text{ tickets in box} - \# \text{ draws}}{\# \text{ tickets} - 1}}$$

Sample of 10 from population of 202

$$\text{correction factor} = \sqrt{\frac{202 - 10}{202 - 1}} = 0.98$$

For a sample of 40 from a population of 202

$$\text{correction factor} = \sqrt{\frac{202 - 40}{202 - 1}} = 0.9$$

ie the chance error should be reduced from the value found earlier (7.7%) and should be

$$7.7 \times 0.9 = ~~6.93~~ \underline{\underline{6.9\%}}$$

When sampling without replacement, the chance variability in % will be ~~big~~ smaller than the chance variability when sampling with replacement.

If we fix the sample size, how big must the population be before the correction is negligible?

Opinion poll of with sample size  $n = 2500$

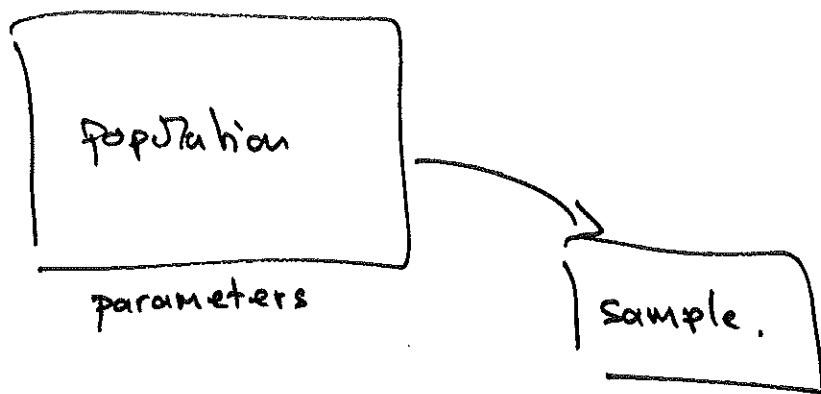
<u>population size</u> (50:50 split)	<u>correction factor</u>
50,000	0.707
10,000	0.866
100,000	0.987
500,000	0.997

← really makes no difference.

So far, we've assumed we've known what was in the box (population) and looked at the characteristics of a sample.

Now: turn it around - given a sample what can we say about the population?

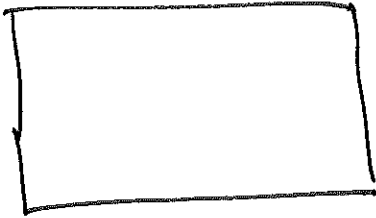
This is called inference.



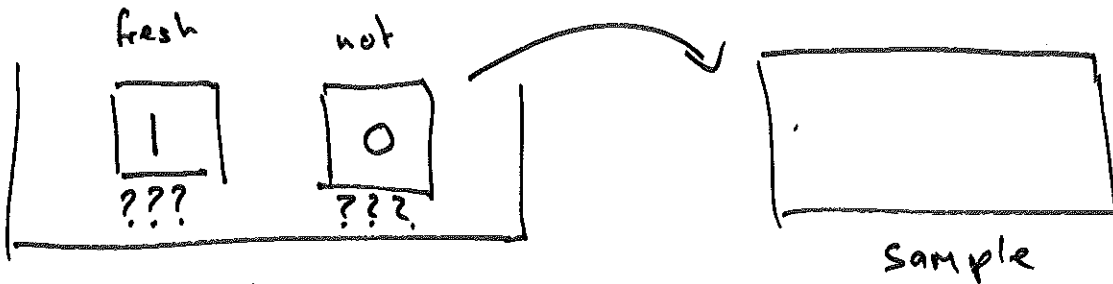
Assume SRS

from the size and composition of the sample, we can say how accurate the estimate of the population parameter will be.

Population: all UCSE students



want to know % fresh.



↑  
 we don't know how many tickets of each type there are (this is what we want to find out!)

Assume we sample 250 students

in our sample we had 70 fresh.

Think about SE. (a measure of how much we would expect the # of fresh in our sample to vary as we take repeated samples)

$$SE = \sqrt{\# \text{ draws}} \times SD_{\text{box}}$$

$$SD_{\text{box}} = (1 - 0) \sqrt{\text{fraction of tickets with 1} \times \text{fraction of tickets with 0}}$$

↖
↗  
 unknown



How to proceed?

- use the fractions in the sample.

$$SD_{\text{box}} = (1-0) \sqrt{\frac{70}{250} \times \frac{180}{250}} = 0.449$$

$$SE \# \text{ fresh.} = \sqrt{250} \times 0.449 = 7.1$$

$$\sqrt{\text{sample size}} \times SD_{\text{box}}$$

$$SE \% \text{ fresh} = \frac{7.1}{250} \times 100 = 2.8\%$$

sample size

$$\begin{aligned} \text{Estimate of the population \%} &\hat{=} \text{sample \%} \\ &= \frac{70}{250} \times 100 = 28\% \end{aligned}$$

So, from our sample of 250 students  
we estimate

% of fresh people on campus is  $28 \pm 2.8\%$ .

- need to be a little bit  
careful in how we interpret  
this.

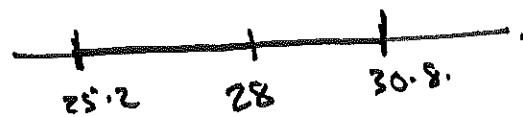
<u>Sample %</u>	=	<u>population %</u>	+	<u>chance error %</u>
28		?		?
28		28		0
28		26		2
28		24		4
28		32		-4

What can we say about population % from a sample % of 28?

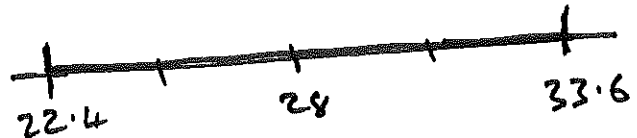
For each value of population % listed, we can think about the resulting chance error in terms of how many SE it is (ie standard units).

We can define a confidence interval around the sample % and use the normal approximation to define the confidence level

sample %  $\pm 1 SE$  is 68% confidence interval



sample %  $\pm 2 SE$  is 95% confidence interval



sample %  $\pm 3 SE$  is 99.7% confidence interval.

" We can be about 95% confident that the % of fresh persons on campus is between 22.4 and 33.6% "

Note: we have not said that the chance of the % of fresh being between 22.4% and 33.6% is 0.95

because: probability is defined as the frequency of occurrence of an event.

the population % is either between 22.4 and 33.6% or it is not.

this does not change however many times you measure it.

The chance is in the sampling, not in the parameter.

the parameter takes a fixed value.

$$\frac{4506}{15,125} \rightarrow 29.8\%$$

From our sample, 70 out of 250 were fresh.

95% CI for population % of 22.4 - 33.6%

In this case, the 95% CI covers the population %.

Consider a 2<sup>nd</sup> sample.

size 250

# fresh 62

$$\text{sample \%} = \frac{62}{250} \times 100 = 24.8\%$$

$$SD_{\text{box}} = (1 - 0) \sqrt{\frac{62}{250} \times \frac{188}{250}} = 0.43$$

↑  
use sample proportion  
in place of unknown  
population proportion

$$\text{SE on Sum} \quad (\# \text{ fresh}) = \sqrt{250} \times 0.43$$

$$\text{SE on \% fresh} = \frac{\sqrt{250} \times 0.43}{250} \times 100 = 2.7\%$$

$$95\% \text{ CI is } 24.8\% \pm 2 \times 2.7\%$$

$$\text{or } 19.4 \rightarrow 30.2\%$$

again, 29.8 is in this range, so this covers  
the population %.

---

Now, repeat this 202 times...

each sample will give a 95% CI.

→ 95% of the CIs cover the true value.

And that is what a 95% CI is:

95% of all samples will generate a CI that  
covers the true value

(but: with just one sample, you don't know  
if your sample does or does not).

