

- Confidence Intervals in General
- Confidence Intervals for Averages
- The Current Population Survey
- Answer Questions

## 16.0 More on Confidence Intervals

This Survey is probably more accurate than a full count would be.

The Bureau of Labor Statistics administers the Current Population Survey (CPS), which is performed by the Census Bureau. The primary purpose of the CPS is to estimate the unemployment rate.

One persistent problem is how to define the category "unemployed". Between 1984 and 1987, England changed their definition about 80 times, and ultimately classified their definition.

## 16.1 The Current Population Survey

The CPS is usually redesigned after every census to achieve greater or comparable accuracy with equivalent or less cost.

One PSU is drawn from each stratum, at random with probability proportional to its population. Some USUs are picked at random from each PSU.

- Counties and cities are grouped into Primary Sampling Units (PSUs).
  - PSUs are grouped into strata according to demographic criteria; each stratum is fairly homogeneous, and does not cross state lines.
  - Each PSU is divided in Ultimate Sampling Units (USUs) of about four housing units.
  - So strata contain PSUs which contain USUs.
- The CPS has the following structure:

The CPS gets about 1 in 1,800 people, with oversampling in small states. The CPS is **not** a simple random sample. The final estimates must be weighted to reflect the unequal probabilities of selection. True or False: The CPS uses Multistage Cluster Sampling. A big concern is the possibility that, during the years between censuses, there might be a significant relocation of population. This happened in the 1980s, when the people left the "Rust Belt" to find employment in the "Sun Belt".

The CPS also uses a **rotating panel**. The same households are visited in four consecutive months, then dropped from the sample for eight months, then revisited for four consecutive months.

However, it also allows **panel bias**. People's answers change systematically over the course of their involvement in a panel. For the CPS, people are more likely to report that they are looking for work on the first response than the second.

- builds higher response rates through steady contact.
- reduces costs—fewer new addresses are drawn
- reduces burnout rates from respondents
- reduces variance by controlling for family effects
- longitudinal study of employment changes

A rotating panel enables

The CPS also uses a **rotating panel**. The same households are visited in four consecutive months, then dropped from the sample for eight months, then revisited for four consecutive months.

These definitions may not adequately capture what people imagine the unemployment rate to truly be. But it is important not to change definitions too casually, since one cannot compare employment trends over time.

someone who has stopped looking for work.

- outside the labor force, meaning a student, a retired person, or for work and had looked within the last four weeks;
- unemployed, meaning the person was not employed but was available was temporarily absent from their job;
- employed, meaning the person did paid work in the previous week or

The basic definitions are as follows:

But this does not control for any biases that may be present.

independent halves.

at the size of the differences between estimates made from those two households in each USU into two groups of two households, and look would arise between two independent studies. One can divide the four method. The standard error in an estimate is the typical difference that To estimate the standard errors, the Census Bureau uses the half-sample

that subcategory becomes quite small.

in those estimates becomes large because the number of respondents in But if the categories get too refined, then the uncertainty (or variance) proportion of unemployed white females, and how this changes over time. The results are broken out into categories. One can estimate the

- Why don't we use a finite population correction factor?
- Recall the standard error of a proportion:  $\sqrt{\frac{p(1-p)}{n}}$ . Does this apply?
- Why don't we use a simple random sample?

Calculations from demographic trends show that the CPS estimate of the size of labor force has standard error about 5% smaller than one would get from a comparable simple random sample. So the weighting helps. But the standard error in the unemployment percentage is about 50% larger than one would get from a comparable simple random sample. So the clustering hurts.

To be a good approximation, the CLT requires that  $n$  be fairly large (about 25 is usually adequate for most purposes).

where  $EV$  is the mean of the box,  $sd$  is the standard deviation of the box, and  $n$  is the number of draws.

$$\underline{X} \sim N(0, 1) \quad \frac{\sqrt{n}/ps}{EV - X}$$

Recall the Central Limit Theorem for averages:

## 16.2 Confidence Intervals for Averages

These standard errors are all very similar ideas, but they look different. Probably it is a good idea to memorize the formulas for the sum, average, and proportion. In a higher-level class, the connections between these and proportion become clear.

- percentage
- proportion
- count
- average
- sum

What are the standard errors of a:

the standard error of the average?

What is the difference between the standard deviation of the sample and

Recall that a **confidence interval** is a random interval  $[L, U]$  such that  $C\%$  of the time, the population average or proportion will be greater than  $L$  but less than  $U$ .

The numbers  $L$  and  $U$  are obtained from the sample by using the CLT.

The analyst gets to pick the **confidence level**  $C$ .

The formula for a confidence interval on a population mean is:

$$\bar{X} \pm z_C * s_e$$

where  $z_C$  is the value from a normal table such that the area between  $-z_C$  and  $-z_C$  is  $C$ .

intervals contain the population mean. which is the case). All we can say is that C% of similarly constructed interval will either contain the true value or not (and we won't know contain the true population mean is C/100. **After** collecting the data, the **Before** setting the CI, one can say that the probability that it will

then we can estimate it by the standard deviation of the sample. If we do not know the standard deviation of the box (or the population),

more quickly. Since  $se = sd/\sqrt{n}$ , the width  $U - L$  of the confidence interval goes to zero as  $n$  increases. If we have sampled without replacement from a finite population, then  $se = FPCF * sd/\sqrt{n}$  and the width goes to zero even zero as  $n$  increases. If we have sampled without replacement from a finite

- the  $sd$  of the population is known.
- when the sample size is sufficiently large, and the (rare) case when enable a good estimate of the  $sd$  of the box (population), the case
- distinguishes the case when the sample size is not sufficiently large to
  - set CIs on proportions and averages
  - set one-sided confidence intervals

Note that one can:

We shall go over the handout entitled “Blueprint for Confidence Intervals”.

## 16.3 Confidence Intervals in General