

- Confidence Intervals
- More on the CLT
- Answer Questions

15.0 The CLT and Confidence Intervals

$$\cdot \sim N(0, 1) . \frac{\underline{n} * ps}{\text{sum} - n * EV}$$

Similarly, for sums, just multiply the left-hand side by n/n to get:

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

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

Recall the Central Limit Theorem for averages:

15.1 More on the CLT

histogram of them, what will the histogram of the averages look like?

If I take 1000 samples of size 20, find the 1000 separate \bar{X} 's, and make a

non-normal distribution.

Suppose the population (i.e., the tickets in the box) has some weird,

$$\cdot \underline{(d - 1)d} \wedge = ps$$

this box of zeroes and ones is:

One can show (in a more advanced class) that the standard deviation for

categories.)

a proportion—the zeroes and ones are just codes for two different and Republicans, of course, since our real interest is in estimating p. (We could replace zeroes and ones by Heads and Tails or Democrats The expected value for the box is the population proportion of ones, or

contains B tickets, labelled as zeroes and ones. The CLT still applies. As an extreme example of a non-normal box, suppose the box just

Suppose one draws a sample of 100 people at random from the U.S. and asks them whether they believe that Karl Rove is a jackal. And suppose 22 of them say yes. What is the probability that the true proportion of Americans who think Rove is a jackal is 25%?

We can use the CLT—the box model has ones for people who hate Rove, zeroes for those who do not. The expected value of the box is

$$EV = 1/B \{ \text{sum of zeroes and ones for all the U.S.} \} = p$$

which is the proportion we are trying to estimate.

Also, our estimate is just a sample average:

$$\hat{p} = 1/100 \sum_{i=1}^{100} X^i = 22/100 = .22$$

where X^i is 1 if the i th respondent hates Rove, and is zero else. Thus the CLT for averages applies.

Note that we used the sample value \hat{p} to estimate the sd of the box.

From the standard normal table, we know this has chance $(1/2)(100 - 51.61) = 24.195\%$, so the probability that the true proportion is more than .25 is about .24.

$$\begin{aligned}
 P[EV < .25] &= P[X - EV < \frac{se}{.25}] \\
 &= P[Z < \frac{\sqrt{p(1-p)/n}}{.25} - \frac{EV - \hat{p}}{\sqrt{p(1-p)/n}}] \\
 &= P[Z < \frac{\sqrt{.22 * (.78)/100}}{.25} - \frac{.22 - .25}{\sqrt{.22 * (.78)/100}}] \\
 &= P[Z < -.724]
 \end{aligned}$$

number of draws.
where B is (as always) the number of tickets in the box and n is the

$$FPCF = \sqrt{\frac{B-1}{B-n}}$$

The finite population correction factor is:

correction factor or FPCF.

(why?), and we should adjust for this by using the finite population replacement shrinks the standard error for the average or the sum replacement without

This can have a big effect in small populations. Sampling

sampled once.

situations, we do not draw with replacement. Respondents are only (or, equivalently, that the population is infinite). But in most survey

Our CLT assumes that one draws from the box with replacement

$$FPCF = \sqrt{\frac{200 - 1}{200}} = .70888.$$

But if we had drawn the random sample from a town of 200, then

$$FPCF = \sqrt{\frac{290,000 - 1}{290,000}} = .9999.$$

the U.S. population. In that case,

For example, in the Karl Rove example, we chose a random sample from

fraction of B , and not much otherwise.

This FPCF will reduce the standard error by a lot when n is a significant

which is .15865, a bit smaller than before.

$$\begin{aligned}
 P[EV < .25] &= P\left[\frac{\text{FPCF} * \sqrt{p(1-p)/n}}{p - .25} > \frac{\text{FPCF} * \sqrt{p(1-p)/n}}{p - .25}\right] \\
 &= P[Z < \frac{\text{FPCF} * \sqrt{p(1-p)/n}}{p - .25}] \\
 &= P[Z < (.22 - .25) / (.70888 * \sqrt{.22 * (1 - .22) / 100})] \\
 &= P[Z < -1.024]
 \end{aligned}$$

For the Karl Rove example in a town of 200,

$$\left. \begin{array}{l} \text{for averages } \frac{sd/\sqrt{n}}{sd * \sqrt{n}} \\ \text{for sums } \end{array} \right\} = es$$

Whenever one samples without replacement from a small population, one should multiply the usual standard error by the FPCF. Recall that:

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

The analyst gets to pick the confidence level C .

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

15.2 Confidence Intervals

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.

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

Since $se = sd/\sqrt{n}$, the width $U - L$ of the confidence interval goes to

and $-z_C$ is C .

where z_C is the value from a normal table such that the area between z_C

$$\text{CZ} * se \pm z_C = U, L$$

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

knew the true sd , then we could solve to find p .)
 $p(1 - p)$, but p is what we are trying to estimate in the first place. If we always have to do this for proportions: the true standard deviation is which is just the standard deviation of the sample. (In general we shall population), then we can estimate the standard deviation by $\hat{p}(1 - \hat{p})$, As before, if we do not know the standard deviation of the box (i.e., the

a bit more readable.

This is the same formula as before, since the sample proportion is just a sample average. All we have done is re-write the formula in a way that is

$$\frac{u}{\sqrt{\hat{p}(1 - \hat{p})}} \pm z^* \cdot \sigma_z$$

Similarly, the formula for a confidence interval on a proportion is:

$$\hat{p} = 18/100 = .18.$$

Your estimate of the proportion of people who have seen the show is

worry about the FPCF?)

You sample 100 people at random; 82 are “virgins.” (Do you need to

adults who have seen “The Rocky Horror Picture Show.”

Suppose you want a 95% confidence interval on the proportion of U.S.

15.3 Problems

Thus the 95% confidence interval on the proportion who have seen the

Similarly, L is found by subtracting $se * z_C$ from \hat{p} , and is .1051.

$$\begin{aligned} & = .2549. \\ & = .18 + \sqrt{\frac{(.18)(.82)}{100}} * (1.95) \\ & = \underline{.02 * \frac{u}{(d - 1)d}} + \underline{d} = U \end{aligned}$$

For $C = 95$, the normal table gives $z_C = 1.95$. So

$$\underline{.02 * \frac{u}{(d - 1)d}} \mp \underline{d} = U$$

A $C\%$ CI on the true p is given by

show is [.1051, .2549].

The reason for this is that the true proportion is either within the interval or it isn't—there is no randomness in the parameter (unless you are a Bayesian...). Instead, the randomness comes from the sample.

So all we can say is that 95% of the time, we will draw a sample that generates a confidence interval that contains the true value.

Instead, one should say that "In 95% of similarly constructed intervals, the true proportion will be within the interval."

One has to be careful when interpreting this confidence interval. It is technically **wrong** to say that the probability is .95 that the true proportion of people who have seen the "Rocky Horror Picture Show" is between .1051 and .2549.

Instead, one should say that "In 95% of similarly constructed intervals,