

Main idea today, part 1:

Means of repeated samples center around the population mean

Examples in Agresti & Finlay reading

Population-of-three example last time

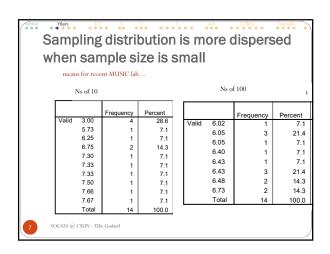
Lab9 & HW3, across the class

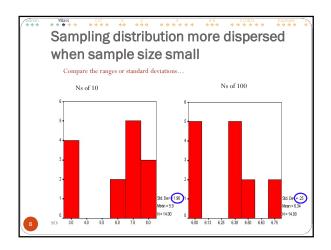
Two illustrations in lecture next time

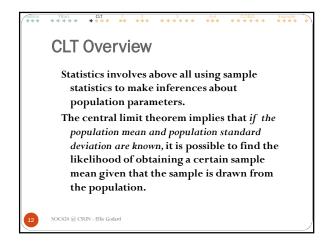
Outline for Today...

Basic Idea of the CLT
Six Implications
Sampling Distributions & Z-scores

An example
A lab
Brilliance & wit interspersed throughout







Eventually, it would be normal

•That's only 14 samples.

•With 1606 GSS cases, there are (1606!/1506!) possible samples of n=100, =(1606x1605x1604...)/(1506x1505x...), a number with over 312 digits!

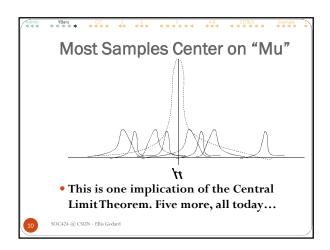
•The more samples we add to that histogram the closer to normal it would look

•And its mean would be .... what? ☺

CLT (as stated in Agresti & Finlay)

Consider a random sample of n measurements from a population distribution having mean  $\mu$  and standard deviation  $\sigma$ .

Then, if n is sufficiently large, the distribution of sample means (i.e. the sampling distribution) is an approximately normal distribution with mean  $\mu$  and standard error  $\sigma_{\overline{\nu}} = \sigma / \sqrt{n}$ 



CLT (in a few more words)

If a sample is taken randomly and if it is sufficiently large, the sampling distribution is normal, even if the data distribution (that is, the sample or population data) is not normal, i.e.

• An estimated 68.26% of sample means fall within 1 standard deviation of the sampling distribution – that is, within one standard error of the population mean

• An estimated 95.44% fall within 2 standard errors

• An estimated 95% fall within 1.96 standard errors

Note that these are estimates because we use the sample standard deviation when we don't know the population parameter (which is often)

#### 3 Main Implications of Theorem:

- sampling distribution's mean is approximately equal to the pop. mean
  - Samples cluster around pop. Mean
  - That's above, and in last 2 set of notes
- Standard error decreases as the sample size increases
  - By formula, since n is denominator (see last 2 lects)
- By logic: larger samples are closer to the pop. mean, but smaller samples could be weirder and further from it
- Regardless of the shape of population distribution, the sampling distribution will be normal

SOC424 @ CSUN - Ellis Godard

## 2<sup>nd</sup> Implication: n versus se

- As n increases, the standard error decreases. Since the square root of n is in the denominator, a larger sample size will decrease the standard error
- See figure 4.14 on p.103 in A&F

SOC424 @ CSUN - Ellis Godard

1st Implication:  $\overline{Y}_{\overline{v}} = \mu$ 

- Imagine two simulations, repeatedly drawing samples from a larger set of cases:
  - I drew 200 samples and got a mean of 500.4
  - I drew 2000 samples and got mean of 499.9.
  - Both are close to the population mean of 500.
- The more cases there are, the closer the samples are on average to the pop parameter
  - As the number of samples increased, we would expect the mean of the sampling distribution to become closer and closer to 500.

SOC424 @ CSUN - Ellis Godard

From above example, if ...

$$n=10$$
, then  $\sigma_{\bar{y}} = \frac{100}{\sqrt{10}} = 31.6$ 

$$n=100$$
, then  $\sigma_{\bar{y}} = \frac{100}{\sqrt{100}} = 10$ 

$$n=1000$$
, then  $\sigma_{\bar{y}} = \frac{100}{\sqrt{1000}} = 3.16$ 

$$n=10000$$
, then  $\sigma_{\bar{y}} = \frac{100}{\sqrt{10000}} = 1$ 

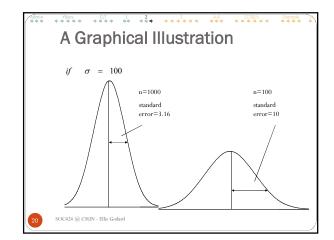
SOC424 @ CSUN - Ellis Godard

Example if  $\sigma = 100$  and n = 1000, then  $\sigma_{\bar{y}} = \sqrt[\sigma]{n}$ 

 $\sigma_{\bar{y}} = \frac{100}{\sqrt{1000}} = \frac{100}{31.6228} = 316$ 

In the above simulations (see last lecture), the computed standard deviations were 3.13 and 3.22. Pretty close!

SOC424 @ CSUN - Ellis Godard



## 3<sup>rd</sup> Implication: Normality

- The fact that the sampling distribution is approximately normal has nothing to do with the shape of the population distribution
- Sampling distribution is always normal
- See Figure 4.12 on page 101 of A&F

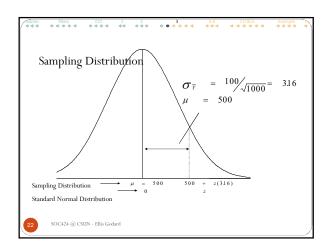
SOC424 @ CSUN - Ellis Godard

### RE II: Political Sociology

- What is the expected shape of the distribution of "presidential support" in a population of highly partisan Congress?
- What is the expected shape of the sampling distribution of sample means if the samples are "sufficiently large"?
  - Note: Next week, we'll define "sufficiently"

24

SOC424 @ CSUN - Ellis Godard



Affinin Vibra CIT 1 2 3 4-6 CITEZE Fample >

### RE III: Status and Orgs

- What is the expected distribution of "academic reputation" for a population of 200 national universities?
- If we selected many samples of sufficient size (25 to 30 is usually enough according to Agresti and Finlay) and computed mean academic reputation for each sample, what is the shape of the distribution of these means? Approximately what is the mean?

25

SOC424 @ CSUN - Ellis Godard

# Rhetorical Example I: Sports

- What is the *expected* distribution of "number of years in the major league" for a baseball player?
- What is the shape of the sampling distribution?

SOC424 @ CSUN - Ellis Godard

## 3 Main Implications of CLT:

- sampling distribution's mean approximately equal to population mean
  - Samples cluster around pop. mean
- standard error decreases as the sample size increases
  - Larger samples closer to pop. mean
- regardless of shape of population distribution, sampling distribution normal



SOC424 @ CSUN - Ellis Godard



6th Implication: Size Matters

If the sample size n is small (less than 25-30), the shape of the sampling distribution is approximately normal only if the population distribution itself is approximately normal

See Figure 4.15 on p. 104 of A&F

30

SOC424 @ CSUN - Ellis Godard

4th Implication: Extreme Cases

1. What is the distribution of the sample mean if n = 1?

The probability distribution for one randomly selected observation is the same as the population distribution of the variable.

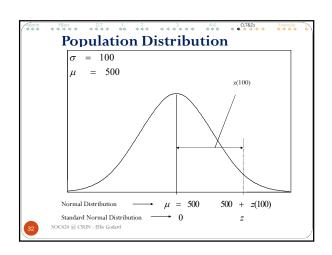
2. What is the distribution of the sample mean if n=N (sample size equal population size)?

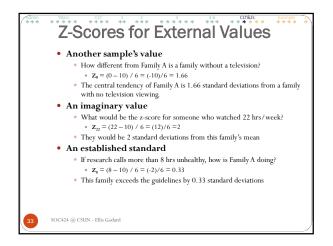
The sample means always equals the population mean; the distribution is concentrated at the population mean μ

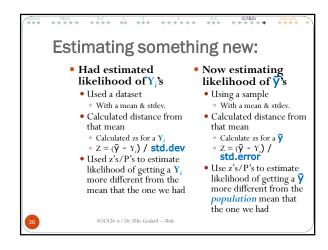
Recall that a z-score corresponding to a particular value is the number of standard deviations that this value is from the mean of the distribution.  $z = \frac{Y - \mu}{\sigma}$ 

5<sup>th</sup> Implication: Parameters Matter

• The larger the standard deviation of a variable in the population, the larger the standard error of the sampling distribution. This again follows from the definition of a standard error:  $\sigma_{\overline{y}} = \sigma / \sqrt{n}$ 







Distance by Standard Errors

• I've illustrated z-scores w/ comparisons we might make (e.g. to an imaginary value, an external standard, or another sample's mean)

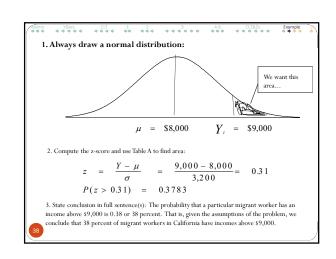
• But we make such comparisons in terms not of standard deviations but of standard errors:  $stderror = \frac{stdev}{\sqrt{n}}$ • For example, for Family A, the standard error is (6 / sqrt of 6), or 2.449

Example from A&F (p89, # 15, 2nd edition)

Suppose the distribution of yearly incomes of all migrant workers in California has a mean of \$8,000 and a standard deviation of \$3,200.

a. Assuming that the distribution is normal, what is the probability that a particular migrant worker has an income over \$9,000?

Specified to a sampling distribution, the z-score of a particular sample mean is the number of standard errors that this sample mean is from the mean of the sampling distribution (which is equal to the population mean).  $z = \frac{\overline{Y} - \mu}{\sigma_{\overline{Y}}}$ 



YBars CLT 1 2 3 4-6 CLT&ZS Exam

### Second part...

(b) If we plan on taking a random sample of 64 migrant workers, what is the sampling distribution of the sample mean income? Find the probability that the sample mean exceeds \$9,000.

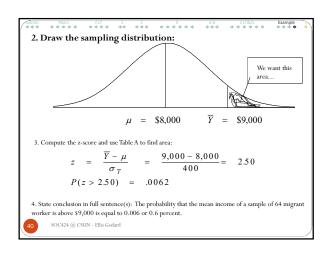
#### **Solution:**

1. Compute the standard error:

$$\sigma_{\overline{y}} = \sqrt[\sigma]{n} = 3,200 / \sqrt{64} = 400$$

39

SOC424 @ CSUN - Ellis Godard



For your next lab...

• Optional extra credit lab

• Count as an additional lab

• Does not replace a skipped/missed lab

• Use the world95.sav dataset

• Note two uses of "population"!!

• For a population of 109 countries

• Not a sample, but all of them

• For each case, its population is measured.

• the number of people in that country

• Group assignment

• 3-5 for secretary bonus

• There's a form on the website.. Use it! ©