

**UCLA STAT 13**  
**Introduction to Statistical Methods for the**  
**Life and Health Sciences**

**Instructor: Ivo Dinov,**  
 Asst. Prof. of Statistics and Neurology

**Teaching Assistants:**  
 Brandi Shanata & Tiffany Head

University of California, Los Angeles, Fall 2007  
[http://www.stat.ucla.edu/~dinov/courses\\_students.html](http://www.stat.ucla.edu/~dinov/courses_students.html)

Slide 1

Stat 13, UCLA, Ivo Dinov

**Chapter 5**  
**Sampling Distributions**

Slide 2

Stat 13, UCLA, Ivo Dinov

**Sampling Distributions**

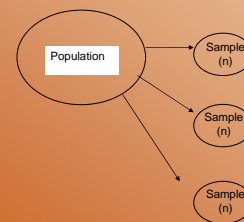
- **Definition:** *Sampling Variability* is the variability among random samples from the same population.
- A probability distribution that characterizes some aspect of sampling variability is called a sampling distribution.
  - tells us how close the resemblance between the sample and the population is likely to be.
- We typically construct a sampling distribution for a statistic.
  - Every statistics has a sampling distribution.

Slide 3

Stat 13, UCLA, Ivo Dinov

**The Meta-Experiment**

- All the possible samples that might be drawn from the population (infinity repetitions).
  - In other words if we were to repeatedly take samples of the same size from the same population, over and over.



Slide 4

Stat 13, UCLA, Ivo Dinov

**The Meta-Experiment**

- Meta-experiments are important because probability can be interpreted as the long run relative frequency of the occurrence of an event.
- Meta-experiments also let us visualize sampling distributions.
  - and therefore understand the variability among the many random samples of a meta-experiment.

Slide 5

Stat 13, UCLA, Ivo Dinov

**Dichotomous Observations**

- Dichotomous - two outcomes
  - (yes or no, good or evil, etc...)
- We use the following notation for a dichotomous outcome
 

$P$	population proportion
$\hat{p}$	sample proportion
- The big question is how close is  $\hat{p}$  to  $P$ ?
- To determine this we need to examine the sampling distribution of  $\hat{p}$
- What we want to know is:
  - if we took many samples of size  $n$  and observed  $\hat{p}$  each time, how would those values of  $\hat{p}$  be distributed around  $p$ ?

Slide 6

Stat 13, UCLA, Ivo Dinov

## Dichotomous Observations

**Example:** Suppose we would like to estimate the true proportion of male students at UCLA. We could take a random sample of 50 students and calculate the sample proportion of males.

- What is the correct notation for:
  - the true proportion of males?
  - the sample proportion of males?
- Suppose we repeat the experiment over and over. Would we get the same proportion of males for the second sample?

Slide 7

Stat 13, UCLA, Jon Dineen

## Reece's Pieces Experiment

**Example:** Suppose we would like to estimate the true proportion of orange reece's pieces in a bag. To investigate we will take a random sample of 10 reece's pieces and count the number of orange. Next we will make an approximation to a sampling distribution with our class results.

What you need to calculate:

- the number of orange
- the sample proportion of orange (number of orange/10)

Slide 8

Stat 13, UCLA, Jon Dineen

## An Application of a Sampling Distribution

**Example:** Mendel's pea experiment. Suppose a tall offspring is the event of interest and that the true proportion of tall peas (based on a 3:1 phenotypic ratio) is  $3/4$  or  $p = 0.75$ . If we were to randomly select samples with  $n = 10$  and  $p = 0.75$  we could create a probability distribution as follows:

	$\hat{p}$	Number Tall	Number Dwarf	Probability
	0.0	0	10	0.000
	0.1	1	9	0.000
	0.2	2	8	0.000
	0.3	3	7	0.003
	0.4	4	6	0.016
	0.5	5	5	0.058
	0.6	6	4	0.146
	0.7	7	3	0.250
	0.8	8	2	0.282
	0.9	9	1	0.188
	1.0	10	0	0.056

Validate using:

[http://www.stat.ucla.edu/Applets/Normal\\_T\\_Ch2\\_F\\_Tables.htm](http://www.stat.ucla.edu/Applets/Normal_T_Ch2_F_Tables.htm)

E.g.,  $B(n=10, p=0.75, a=6, b=6)=0.146$

Slide 9

Stat 13, UCLA, Jon Dineen

## An Application of a Sampling Distribution

- What is the probability that 5 are tall and 5 are dwarf?

$$P(5 \text{ tall and } 5 \text{ dwarf}) = P(\hat{p} = 5/10)$$

$$= P(\hat{p} = 0.5)$$

$$= 0.058$$

	$\hat{p}$	Number Tall	Number Dwarf	Probability
	0.0	0	10	0.000
	0.1	1	9	0.000
	0.2	2	8	0.000
	0.3	3	7	0.003
	0.4	4	6	0.016
→	0.5	5	5	0.058
	0.6	6	4	0.146
	0.7	7	3	0.250
	0.8	8	2	0.282
	0.9	9	1	0.188
	1.0	10	0	0.056

Slide 10

Stat 13, UCLA, Jon Dineen

## An Application of a Sampling Distribution

- If we think about this in terms of a meta-experiment and we sample 10 offspring over and over, about 5.8% of the  $\hat{p}$ 's will be 0.5.

- This is the sampling distribution of sample proportion of tall offspring is the distribution of in repeated samples of size 10.

- If we take a random sample of size 10, what is the probability that six or more offspring are tall?

$$P(\hat{p} \geq 0.6) = 0.146 + 0.250 + 0.282 + 0.188 + 0.056 = 0.922$$

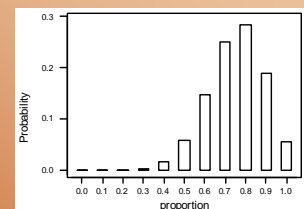
Slide 11

Stat 13, UCLA, Jon Dineen

## An Application of a Sampling Distribution

- This table could also be represented as a histogram with probability on the y-axis and proportion on the x-axis.

- easier to draw these by hand



Slide 12

Stat 13, UCLA, Jon Dineen

## Relationship to Statistical Inference

- We can also use our sampling distribution of  $\hat{p}$  to estimate how much sampling error there is within 5 percentage points of  $p$ . Because we knew  $p$  from the previous example ( $p=0.75$ ), we might want to estimate:

$$P(0.7 \leq \hat{p} \leq 0.8) \\ = 0.250 + 0.282 = 0.532$$

There is a 53% chance that for a sample of size 10,  $\hat{p}$  will be within  $\pm 0.05$  of  $p$ .

- This seems a little crazy, why?

$\hat{p}$	Number Tall	Number Dwarf	Probability
0.0	0	10	0.000
0.1	1	9	0.000
0.2	2	8	0.000
0.3	3	7	0.003
0.4	4	6	0.016
0.5	5	5	0.058
0.6	6	4	0.146
0.7	7	3	0.250
0.8	8	2	0.282
0.9	9	1	0.188
1.0	10	0	0.056

Slide 13

Stat 13, UCLA, Jon Dineen

## Relationship to Statistical Inference

- So far we have been using  $p$  to determine the sampling distribution of  $\hat{p}$ .
- Why sample for  $\hat{p}$  when we already know  $p$ ?
  - We don't need to know  $p$  to get a good estimate (this will come later).

Slide 14

Stat 13, UCLA, Jon Dineen

## Sample Size

- As  $n$  gets larger,  $\hat{p}$  will become a better estimate of  $p$ .
- Just to show...

N	$P(0.7 \leq \hat{p} \leq 0.8)$
10	0.53
20	0.56
50	0.673
100	0.798

\*These calculations were done using the SOCR binomial distribution Calculator.

[http://www.socr.ucla.edu/Applets.dir/Normal\\_T\\_Chi2\\_F\\_Tables.htm](http://www.socr.ucla.edu/Applets.dir/Normal_T_Chi2_F_Tables.htm)

E.g.,  $B(n=20, p=0.75, a=0.7 \times 20=14, b=0.8 \times 20=16)=0.5606$   
**THE POINT:** A larger sample improves the chance that  $\hat{p}$  is close to  $p$ .

- Caution: this doesn't necessarily mean that the estimate will be closer to  $p$ , only that there is a better chance that it will be close to  $p$ .

Slide 15

Stat 13, UCLA, Jon Dineen

## Sampling Distribution for the Mean and Introduction to Confidence Intervals

Slide 16

Stat 13, UCLA, Jon Dineen

## Quantitative Data

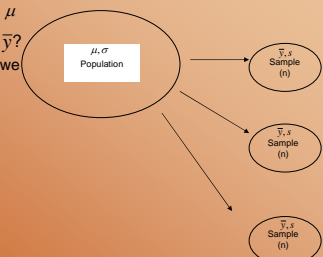
- More complex than dichotomous data
- Sample and populations for quantitative data can be described in various ways: mean, median, standard deviation (each has its own sampling distribution.)

Slide 17

Stat 13, UCLA, Jon Dineen

## Sampling Distribution of $\bar{y}$

- Recall:  $\bar{y}$  is used to estimate  $\mu$
- Question: How close to  $\mu$  is  $\bar{y}$ ?
  - Before we can answer this we need to define the probability distribution that describes sampling variability of  $\bar{y}$



Slide 18

Stat 13, UCLA, Jon Dineen

## Sampling Distribution of $\bar{y}$

- Two really important facts:
  - The average of the sampling distribution of  $\bar{y}$  is  $\mu$ 
    - Notation:  $\mu_{\bar{y}} = \mu$
  - The standard deviation of the sampling distribution of  $\bar{y}$  is  $\frac{\sigma}{\sqrt{n}}$ 
    - Notation:  $\sigma_{\bar{y}} = \frac{\sigma}{\sqrt{n}}$
- Note: As  $n \rightarrow \infty$ ,  $\sigma_{\bar{y}}$  gets smaller
- Why? Look at the formula
- Intuitively does this make sense?

Slide 19

Stat 13, UCLA, Jon Dineen

## Sampling Distribution of $\bar{y}$

- Theorem 5:1 p.159**
  - $\mu_{\bar{y}} = \mu$  (mean of the sampling distribution of  $\bar{y} = \mu$  the population mean)
  - $\sigma_{\bar{y}} = \frac{\sigma}{\sqrt{n}}$  (standard deviation (sd) of the sampling distribution of  $\bar{y} = \frac{\sigma}{\sqrt{n}}$  the population SD divided by  $\sqrt{n}$ )
- Shape:
  - If the distribution of  $Y$  is normal the sampling distribution of  $\bar{y}$  is normal.
  - Central Limit Theorem (CLT) - If  $n$  is large, then the sampling distribution of  $\bar{y}$  is approximately normal, even if the population distribution of  $Y$  is not normal.

Slide 20

Stat 13, UCLA, Jon Dineen

## Central Limit Theorem (CLT)

- No matter what the distribution of  $Y$  is, if  $n$  is large enough the sampling distribution of  $\bar{y}$  will be approximately normally distributed
  - HOW LARGE??? Rule of thumb  $n \geq 30$ .
- The closeness of  $\bar{y}$  to  $\mu$  depends on the sample size
- The more skewed the distribution, the larger  $n$  must be before the normal distribution is an adequate approximation of the shape of sampling distribution of  $\bar{y}$
- Why?

Slide 21

Stat 13, UCLA, Jon Dineen

## Central Limit Theorem – theoretical formulation

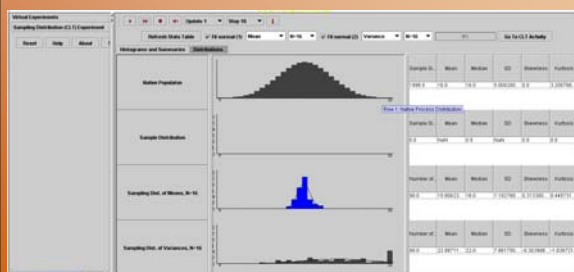
Let  $\{X_1, X_2, \dots, X_k, \dots\}$  be a sequence of **independent** observations from **one specific random process**. Let and  $E(X) = \mu$  and  $SD(X) = \sigma$  and both be finite ( $0 < \sigma < \infty$ ;  $|\mu| < \infty$ ). If  $\bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$  **sample-avg**, Then  $\bar{X}_n$  has a **distribution** which approaches  $N(\mu, \sigma^2/n)$ , as  $n \rightarrow \infty$ .

Slide 22

Stat 13, UCLA, Jon Dineen

## Central Limit Theorem – Empirical validation Sampling Distribution (CLT) Experiment

[http://www.socr.ucla.edu/htmls/SOCR\\_Experiments.html](http://www.socr.ucla.edu/htmls/SOCR_Experiments.html)



Slide 23

Stat 13, UCLA, Jon Dineen

## Linear Combination

Given a collection of  $n$  random variables  $X_1, \dots, X_n$  and  $n$  numerical constants  $a_1, \dots, a_n$ , the rv

$$Y = a_1X_1 + \dots + a_nX_n = \sum_{i=1}^n a_iX_i$$

is called a **linear combination** of the  $X_i$ 's.

Slide 24

Stat 13, UCLA, Jon Dineen

## Expected Value of a Linear Combination

Let  $X_1, \dots, X_n$  have mean values  $\mu_1, \mu_2, \dots, \mu_n$  and variances of  $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$ , respectively

Whether or not the  $X_i$ 's are independent,

$$E(a_1 X_1 + \dots + a_n X_n) = a_1 E(X_1) + \dots + a_n E(X_n) \\ = a_1 \mu_1 + \dots + a_n \mu_n$$

Slide 25

Stat 13, UCLA, Jon Dineen

## Variance of a Linear Combination

If  $X_1, \dots, X_n$  are independent,

$$V(a_1 X_1 + \dots + a_n X_n) = a_1^2 V(X_1) + \dots + a_n^2 V(X_n) \\ = a_1^2 \sigma_1^2 + \dots + a_n^2 \sigma_n^2$$

and

$$\sigma_{a_1 X_1 + \dots + a_n X_n} = \sqrt{a_1^2 \sigma_1^2 + \dots + a_n^2 \sigma_n^2}$$

Slide 26

Stat 13, UCLA, Jon Dineen

## Variance of a Linear Combination

For any  $X_1, \dots, X_n$ , (dependent or independent!!!)

$$V(a_1 X_1 + \dots + a_n X_n) = \sum_{i=1}^n \sum_{j=1}^n a_i a_j \text{Cov}(X_i, X_j)$$

If  $X_i \sim D_i(\mu_i, \sigma_i)$ ,  $X_j \sim D_j(\mu_j, \sigma_j)$  and  $f_{i,j}(x_i, x_j)$

is the joint density function of  $(X_i, X_j)$ , then :

$$\text{Cov}(X_i, X_j) = E((X_i - \mu_i)(X_j - \mu_j)) =$$

$$\int (x_i - \mu_i)(x_j - \mu_j) f_{i,j}(x_i, x_j) dx_i dx_j = E(X_i \times X_j) - \mu_i \mu_j.$$

$$\text{Corr}(X_i, X_j) = \frac{\text{Cov}(X_i, X_j)}{\sigma_i \sigma_j}$$

Slide 27

Stat 13, UCLA, Jon Dineen

## A special case – Difference Between Two Random Variables

If  $X_1 \sim D_1(\mu_1, \sigma_1)$ ,  $X_2 \sim D_2(\mu_2, \sigma_2)$  and  $f_{1,2}(x_1, x_2)$

is the joint density function of  $(X_1, X_2)$ , then :

$$\text{Cov}(X_1, X_2) = E((X_1 - \mu_1)(X_2 - \mu_2)) =$$

$$\int (x_1 - \mu_1)(x_2 - \mu_2) f_{1,2}(x_1, x_2) dx_1 dx_2 = E(X_1 \times X_2) - \mu_1 \mu_2.$$

$$\text{Corr}(X_1, X_2) = \frac{\text{Cov}(X_1, X_2)}{\sigma_1 \sigma_2}$$

Slide 28

Stat 13, UCLA, Jon Dineen

## Example: Variance of Linear Combinations

If  $X_1 \sim T_1(0,1)$ ,  $f_1(x_1) = 2x_1$ ,  $x_1 \in (0,1)$

$X_2 \sim T_2(0,1)$ ,  $f_2(x_2) = -2(x_2 - 1)$ ,  $x_2 \in (0,1)$  and

$f_{1,2}(x_1, x_2) = x_1 + x_2$ ,  $(x_1, x_2) \in (0,1) \times (0,1)$

is the joint density function of  $(X_1, X_2)$ .

$$\text{Cov}(X_1, X_2) = E((X_1 - \mu_1)(X_2 - \mu_2)) =$$

$$\int_0^1 \int_0^1 (x_1 - \mu_1)(x_2 - \mu_2) f_{1,2}(x_1, x_2) dx_1 dx_2 = E(X_1 \times X_2) - \mu_1 \mu_2$$

$$\int_0^1 \int_0^1 (x_1 * x_2) * (x_1 + x_2) dx_1 dx_2 - \left( \int_0^1 x_1 * 2x_1 dx_1 \right) \left( \int_0^1 x_2 * (-2x_2 + 2) dx_2 \right)$$

$$\text{Cov}(X_1, X_2) = \frac{1}{3} - \frac{2}{3} \times \frac{1}{3} = \frac{1}{9}$$

$$\text{Corr}(X_1, X_2) = \frac{\text{Cov}(X_1, X_2)}{\sigma_1 \sigma_2} = \frac{\frac{1}{9}}{\sqrt{\frac{1}{18}} \times \sqrt{\frac{1}{18}}} = 0.5$$

Slide 29

Stat 13, UCLA, Jon Dineen

## Practice Example: Variance of Linear Combinations

If  $X_1 \sim T_1(0,1)$ ,  $f_1(x_1) = x_1$ ,  $x_1 \in (0,2)$

$X_2 \sim T_2(0,1)$ ,  $f_2(x_2) = \frac{3}{2}(x_2 - 1)^2$ ,  $x_2 \in (0,2)$  and

$f_{1,2}(x_1, x_2) = \frac{1}{12}(x_1 x_2 + x_1 + x_2)$ ,  $(x_1, x_2) \in (0,2) \times (0,2)$

is the joint density function of  $(X_1, X_2)$ .

$$\text{Cov}(X_1, X_2) = E((X_1 - \mu_1)(X_2 - \mu_2)) =$$

$$\int_0^2 \int_0^2 (x_1 - \mu_1)(x_2 - \mu_2) f_{1,2}(x_1, x_2) dx_1 dx_2 = E(X_1 \times X_2) - \mu_1 \mu_2$$

$$\int_0^2 \int_0^2 (x_1 * x_2) * \frac{1}{12}(x_1 x_2 + x_1 + x_2) dx_1 dx_2 - \left( \int_0^2 x_1 * x_1 dx_1 \right) \left( \int_0^2 x_2 * \frac{3}{2}(x_2 - 1)^2 dx_2 \right)$$

$$\text{Cov}(X_1, X_2) = ???$$

$$\text{Corr}(X_1, X_2) = \frac{\text{Cov}(X_1, X_2)}{\sigma_1 \sigma_2} = ?$$

Slide 30

Stat 13, UCLA, Jon Dineen



### Application to Data

**Example:** LA freeway commuters (mean/SD systolic pressure):

$$\mu = 130$$

$$\sigma = 20$$

Suppose we randomly sample 4 drivers.

Find  $\mu_{\bar{y}}$   $\mu_{\bar{y}} = \mu = 130$

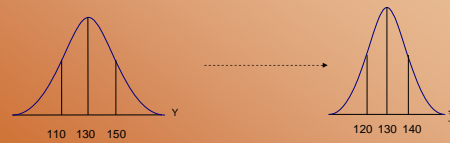
Find  $\sigma_{\bar{y}}$   $\sigma_{\bar{y}} = \frac{\sigma}{\sqrt{n}} = \frac{20}{\sqrt{4}} = 10$

Slide 31

Stat 13, UCLA, Jon Dineen

### Application to Data

● Visually:



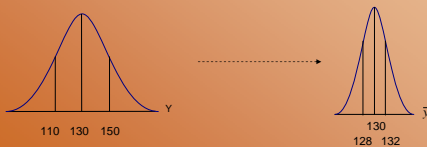
Slide 32

Stat 13, UCLA, Jon Dineen

### Application to Data

**Example:** LA freeway commuters (cont')

Suppose we randomly select 100 drivers



As  $n$  gets larger the variability in the sampling distribution gets smaller.

Slide 33

Stat 13, UCLA, Jon Dineen

### Application to Data

**Example:** LA freeway commuters (cont')

Suppose we want to find the probability that the mean of the 100 randomly selected drivers is more than 135 mmHg

● First step: Rewrite with notation!

$$\bar{y} \sim N(130, 2)$$

● Second step: Identify what we are trying to solve!

$$P(\bar{y} > 135)$$

Slide 34

Stat 13, UCLA, Jon Dineen

### Application to Data

● Third step: Standardize

$$P(\bar{y} > 135) = P\left(\frac{\bar{y} - \mu_{\bar{y}}}{\sigma_{\bar{y}}} > \frac{135 - 130}{2}\right) = P(Z > 2.5)$$

● Fourth Step: Use the standard normal table to solve  
 $1 - 0.9938 = 0.0062$

If we were to choose many random samples of size 100 from the population about 0.6% would have a mean SBP more than 135 mmHg.

Slide 35

Stat 13, UCLA, Jon Dineen

### Application to Data

**Example:** LA freeway commuters (cont')

n	$P(125 < \bar{Y} < 135)$	$\sigma_{\bar{y}}$
4	$P(-0.5 < Z < 0.5) = 0.3830$	$\frac{20}{\sqrt{4}} = 10$
10	$P(-0.79 < Z < 0.79) = 0.5704$	$\frac{20}{\sqrt{10}} = 6.32$
20	$P(-1.12 < Z < 1.12) = 0.7372$	$\frac{20}{\sqrt{20}} = 4.47$
50	$P(-1.77 < Z < 1.77) = 0.9232$	$\frac{20}{\sqrt{50}} = 2.83$

The mean of a larger sample is not necessarily closer to  $\mu$ , than the mean of a smaller sample, but it has a greater probability of being closer to  $\mu$ .

Therefore, a larger sample provides more information about the population mean

Slide 36

Stat 13, UCLA, Jon Dineen

## Notation

- Notation:

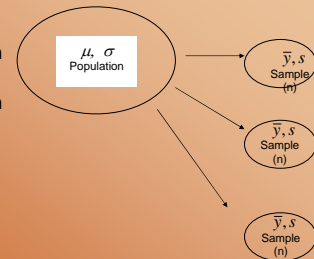
	mean	standard deviation
Population	$\mu$	$\sigma$
Sample	$\bar{y}$	$s$
Sampling Distribution of $\bar{y}$	$\mu_{\bar{y}}$	$\sigma_{\bar{y}} = \frac{\sigma}{\sqrt{n}}$

Slide 37

Stat 13, UCLA, Jon Dineen

## Other Aspects of Sampling Variability

- Sampling variability in the shape
- Sampling variability in the sample standard deviation



Overall: If  $n$  is large,  $s \rightarrow \sigma$ , the shape of each sample will be close to the shape of the population, and the shape of the sampling distribution of  $\bar{y}$  will approach a normal distribution.

Slide 38

Stat 13, UCLA, Jon Dineen

## Statistical Estimation

- This will be our first look at statistical inference
- Statistical estimation is a form of statistical inference in which we use the data to:
  - determine an estimate of some feature of the population
  - assess the precision of the estimate

Slide 39

Stat 13, UCLA, Jon Dineen

## Statistical Estimation

**Example:** A random sample of 45 residents in LA was selected and IQ was determined for each one. Suppose the sample average was 110 and the sample standard deviation was 10.

What do we know from this information?

$$\bar{y} = 110$$

$$S = 10$$

Slide 40

Stat 13, UCLA, Jon Dineen

## Statistical Estimation

- The population IQ of LA residents could be described by  $\mu$  and  $\sigma$
- 110 is an estimate of  $\mu$
- 10 is an estimate of  $\sigma$
- We know there will be some sampling error affecting our estimates
  - Not necessarily in the measurement of IQ, but because only 45 residents were sampled

Slide 41

Stat 13, UCLA, Jon Dineen

## Statistical Estimation

- QUESTION: How good is  $\bar{y}$  as an estimate of  $\mu$ ?
- To answer this we need to assess the reliability of our estimate  $\bar{y}$
- We will focus on the behavior of  $\bar{y}$  in repeated sampling
  - Our good friend, the sampling distribution of  $\bar{y}$

Slide 42

Stat 13, UCLA, Jon Dineen

### The Standard Error of the Mean

- We know the discrepancy between  $\mu$  and  $\bar{y}$  from sampling error can be described by the sampling distribution of  $\bar{y}$ , which uses  $\sigma_{\bar{y}}$  to measure the variability

■ Recall:  $\sigma_{\bar{y}} = \frac{\sigma}{\sqrt{n}}$

- Is there a problem with obtaining  $\sigma_{\bar{y}}$  from our data?
- What seems like a good estimate for  $\sigma_{\bar{y}}$ ?

$\frac{s}{\sqrt{n}}$  is an estimate for  $\frac{\sigma}{\sqrt{n}}$

Called the standard error of the mean

Slide 43

Stat 13, UCLA, Jon Dineen

### The Standard Error of the Mean

- Notation for the standard error of the mean

$$SE_{\bar{y}} = \frac{s}{\sqrt{n}}$$

- Sometimes referred to as the standard error (SE)
- Round to two significant digits

Slide 44

Stat 13, UCLA, Jon Dineen

### The Standard Error of the Mean

**Example:** LA IQ (cont')

$$SE_{\bar{y}} = \frac{10}{\sqrt{45}} = 1.49$$

- What does this mean?
  - Because the standard error is an estimate of  $\sigma_{\bar{y}}$ , it is a measure of reliability of  $\bar{y}$  as an estimate of  $\mu$ .
  - We expect  $\bar{y}$  to be within one SE of  $\mu$  most of the time

Slide 45

Stat 13, UCLA, Jon Dineen

### The Standard Error of the Mean

- If SE is small we have a more precise estimate
- The formula for SE uses  $s$  (a measure of variability) and  $n$  (the sample size)
  - Both affect reliability.

**Example:** LA IQ (cont')

$s$  describes variability from one person in the sample to the next SE describes variability associated with the mean (our measure of precision for the estimate)

Slide 46

Stat 13, UCLA, Jon Dineen

### The Standard Error of the Mean

	n	$\bar{y}$	SE	s
Female	5	117	6.40	14.3
Male	40	109	3.16	20.0

As  $n \rightarrow \infty$ ,

$\bar{y} \rightarrow \mu$

$s \rightarrow \sigma$

$SE \rightarrow 0$

- Example: LA IQ (cont')

Suppose the results of the

45 LA residents were analyzed by gender.

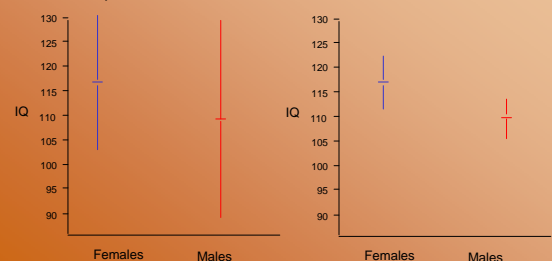
Females have greater variability, but a much smaller SE because their sample size is larger. Therefore the females will have a more reliable estimate of  $\mu$ .

Slide 47

Stat 13, UCLA, Jon Dineen

### The Standard Error of the Mean

- Which plot represents the SD and the SE?
- Which plot describes the data better?



Slide 48

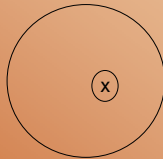
Stat 13, UCLA, Jon Dineen



### Confidence Interval for $\mu$

**Example:** (Analogy from Cartoon Guide to Statistics)  
Consider an archer shooting at a target. Suppose she hits the bulls eye (a 10 cm radius) 95% of the time. In other words, she misses the bulls eye one out of 20 arrows. Sitting behind the target is another person who can't see the bull's eye. The archer shoots a single arrow and it lands:

The person behind the target circles the arrow with a 10 cm radius circle, reasoning that with the archers 95% hit rate, the true center of bull's eye should be within part of that circle.

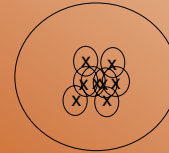


Slide 49

Stat 13, UCLA, Jon Dineen

### Confidence Interval for $\mu$

As she shoots more and more arrows, the person draws more and more circles and finally reasons that these circles will include the true center of the bull's eye 95% of the time.



Slide 50

Stat 13, UCLA, Jon Dineen

### Confidence Interval for $\mu$

- Basic idea of a confidence interval:
  - $\mu$  is the true center of the bull's eye, but we don't actually know where it is
  - We do know  $\bar{y}$ , which is where the arrow came through
  - We can use  $\bar{y}$  and SE from the data to construct an interval that we hope will include  $\mu$

Slide 51

Stat 13, UCLA, Jon Dineen

### Confidence Interval for $\mu$

- Let's build this interval
  - From the standard normal distribution we know:  
 $P(-1.96 < Z < 1.96) = 0.95$
  - How can we rearrange this interval so that  $\mu$  is in the middle?
- Proof
- Formula  $\bar{y} \pm 1.96 \left( \frac{\sigma}{\sqrt{n}} \right)$ 
  - will contain  $\mu$  for 95% of all samples
- Any problems with using this formula with our data?
  - We can use  $s$  to estimate  $\sigma$ , but this changes things a little bit

Slide 52

Stat 13, UCLA, Jon Dineen

### The T Distribution

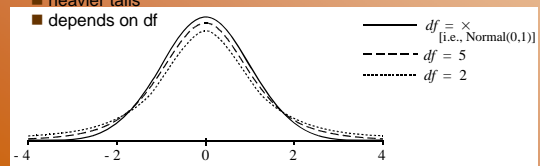
- If the data came from a normal population and we replace  $\sigma$  with  $s$ , we only need to change the 1.96 with a suitable quantity  $t_{0.025}$  from the T distribution
  - Student aka William Gosset (early 1900's)
- The T distribution is a continuous distribution which depends on the degrees of freedom ( $df = n-1$ , in this case) because of the replacement we made with  $s$ : Cauchy ( $df=1$ )  $\rightarrow T_{df} \rightarrow N(0,1)$ ,  $df = \infty$

Slide 53

Stat 13, UCLA, Jon Dineen

### The T Distribution

- As  $n$  approaches  $\infty$ , the  $t$  distribution approaches a normal distribution
- Similarities to the normal distribution include:
  - symmetric
  - centered at 0
- Differences from the normal distribution include:
  - heavier tails
  - depends on  $df$



Slide 54

Stat 13, UCLA, Jon Dineen

## The T Table

- Table 4, p. 677 or back cover of book & Online at SOCR
- <http://socr.ucla.edu/Applets.dir/T-table.html>
- [http://socr.ucla.edu/htmls/SOCR\\_Distributions.html](http://socr.ucla.edu/htmls/SOCR_Distributions.html)
- To use the table keep in mind:
  - table works in the upper half of the distribution (above 0)
  - gives you upper tailed areas
    - this means that the "t scores" will always be positive
    - what do you do if you need a lower tail area?
  - depends on df

Slide 55

Stat 13, UCLA, Jon Dineen

## Using The T Table for CI's

- To use the t table for confidence intervals we will be looking up a "t multiplier" for an interval with a certain level, in this example 95%, of confidence
  - notation for a "t multiplier" is  $t(df)_{\alpha/2}$
  - $t_{0.025}$  (aka  $t_{\alpha/2}$ ) is known as "two tailed 5% critical value"
    - the interval between  $-t_{0.025}$  and  $t_{0.025}$ , the area in between totals 95%, with 5% (aka  $\alpha$ ) left in the tails
  - If we look at the table in the back of the book we'll find:
    - $t_{0.025}$  in the 0.025 column
    - two-tailed confidence level of 95% is at the bottom of the 0.025 column
  - This is half the battle, we still need to deal with df!

Slide 56

Stat 13, UCLA, Jon Dineen

## Using The T Table for CI's

**Example:** Suppose we wanted to find the "t multiplier" for a 95% confidence interval with  $df = 12$

$$t(12)_{0.025} = 2.179$$

<http://socr.ucla.edu/Applets.dir/T-table.html>

- Recall: as  $n \rightarrow \infty$  the t distribution approaches the standard normal distribution
  - also df
  - If we look at the bottom of the table when  $df = \infty$ , the t multiplier for a 95% CI is 1.960
    - Does anything seem familiar about this?

Slide 57

Stat 13, UCLA, Jon Dineen

## Calculating a CI for $\mu$

- To calculate a  $100(1 - \alpha)$  CI for  $\mu$ :
  - choose confidence level (for example 95%)
  - take a random sample from the population
    - must be reasonable to assume that the population is normally distributed
  - compute:  $\bar{y} \pm t(df)_{\alpha/2} \left( \frac{s}{\sqrt{n}} \right)$
- Where  $100(1 - \alpha)$  is the desired confidence
  - This means that for a 95% confidence interval  $\alpha$  is 0.05 (or 5%, because  $100(1-0.05) = 0.95$ )

Slide 58

Stat 13, UCLA, Jon Dineen

## Application to Data

**Example:** Suppose a researcher wants to examine CD4 counts for HIV(+) patients seen at his clinic. He randomly selects a sample of  $n = 25$  HIV(+) patients and measures their CD4 levels (cells/uL). Suppose he obtains the following results:

Descriptive Statistics: CD4

Variable	N	N*	Mean	SE Mean	StDev	Minimum	Q1	Median	Q3	Maximum
CD4	25	0	321.4	14.8	73.8	208.0	261.5	325.0	394.0	449.0

Calculate a 95% confidence interval for  $\mu$

Slide 59

Stat 13, UCLA, Jon Dineen

## Application to Data

- What do we know from the background information?

$$\begin{aligned}\bar{y} &= 321.4 \\ s &= 73.8 \\ SE &= 14.8 \\ n &= 25\end{aligned}$$

$$\begin{aligned}\bar{y} \pm t(df)_{\alpha/2} \left( \frac{s}{\sqrt{n}} \right) &= 321.4 \pm t(24)_{0.05/2} \left( \frac{73.8}{\sqrt{25}} \right) \\ &= 321.4 \pm 2.064(14.8) = 321.4 \pm 30.547 \\ &= (290.85, 351.95)\end{aligned}$$

Slide 60

Stat 13, UCLA, Jon Dineen

### Application to Data

- (290.85, 351.95) – great!
- What does this mean?
  - CONCLUSION: We are highly confident at the 0.05 level (95% confidence), that the true mean CD4 level in HIV(+) patients at this clinic is between 278.58 and 342.82 cells/uL.
- Important parts of a CI conclusion:
  1. Confidence level (alpha)
  2. Parameter of interest
  3. Variable of interest
  4. Population under study
  5. Confidence interval with appropriate units

Slide 61

Stat 13, UCLA, Jon Dineen

### Application to Data

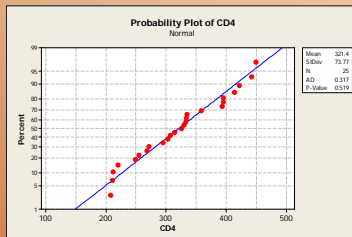
- Still, does this CI (290.85, 351.95) mean anything to us? Consider the following information:
  - The U.S. Government classification of AIDS has three official categories of CD4 counts –
    - asymptomatic = greater than or equal to 500 cells/uL
    - AIDS related complex (ARC) = 200-499 cells/uL
    - AIDS = less than 200 cells/uL
- Now how can we interpret our CI?

Slide 62

Stat 13, UCLA, Jon Dineen

### Application to Data

- Another important point to remember is that our CI was calculated assuming that the data we collected came from a population that was normally distributed!
  - N = 25 so the CLT does not protect us
  - How can we check this?



Slide 63

Stat 13, UCLA, Jon Dineen

### CI Interpretation

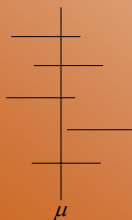
- If we were to perform a meta-experiment, and compute a 95% confidence interval about for each sample, 95% of the confidence intervals would contain  $\mu$
- We hope ours is one of the lucky ones that actually contains  $\mu$ , but never actually know if it does
- We can interpret a confidence interval as a probability statement if we are careful!
  - OK:  $P(\text{the next sample will give a CI that contains } \mu) = 0.95$ 
    - random has happened yet
  - NOT OK:  $P(291 < \mu < 352) = 0.95$ 
    - not random anymore, either  $\mu$  is in there or it isn't

Slide 64

Stat 13, UCLA, Jon Dineen

### CI Interpretation

- The confidence level is a property of the method rather than of a particular interval
- [http://socr.ucla.edu/htmls/SOCR\\_Experiments.html](http://socr.ucla.edu/htmls/SOCR_Experiments.html) → CI



Slide 65

Stat 13, UCLA, Jon Dineen

### Other CI Levels

**Example:** CD4 (cont')

What if we calculate a 90% confidence interval for  $\mu$

- Without recalculating, will this interval be wider or narrower?
- NOTE: Using the same data as before, the only part that changed was the t multiplier.
  - 95%:  $t(24)_{0.025} = 2.064$
  - 90%:  $t(24)_{0.05} = 1.711$
  - As our confidence goes down the interval becomes narrower (because t gets smaller)
  - As the confidence goes up the interval becomes wider

Slide 66

Stat 13, UCLA, Jon Dineen

### Other CI Levels

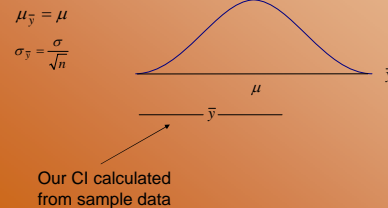
- However, we are sacrificing confidence
  - A 50% CI would be nice and small, but think about the confidence level!
- Better solution: We can also increase the sample size which will make the confidence interval narrower at the same level.
  - Why does this work?

Slide 67

Stat 13, UCLA, Jon Dineen

### Relationship to the Sampling Distribution of $\bar{y}$

- Recall: A CI will contain  $\mu$  for 95% of samples (in repeated sampling, at 95% confidence)



Slide 68

Stat 13, UCLA, Jon Dineen

### Example

**Example:** A biologist obtained body weights of male reindeer from a herd during the seasonal round-up. He measured the weight of a random sample of 102 reindeer in the herd, and found the sample mean and standard deviation to be 54.78 kg and 8.83 kg, respectively. Suppose these data come from a normal distribution.

Calculate a 99% confidence interval.

Slide 69

Stat 13, UCLA, Jon Dineen

### Example

$$\begin{aligned}\bar{y} \pm t(df)_{\alpha/2} \left( \frac{s}{\sqrt{n}} \right) &= 54.78 \pm t(101)_{0.005} \left( \frac{8.83}{\sqrt{102}} \right) \\ &= 54.78 \pm (2.626)(0.874) \\ &= 54.78 \pm 2.296 \\ &= (52.48, 57.08)\end{aligned}$$

- CONCLUSION: We are highly confident, at the 0.01 level, that the true mean weight of male reindeer from the herd during this seasonal round-up is between 52.48 and 57.08 kg.

Slide 70

Stat 13, UCLA, Jon Dineen

### Example – 5.39 (in the textbook)

- Suppose proportion of blood type O is 0.44. If we take a random sample of 12 subjects and make a note of their blood types what is the probability that exactly 6 subjects have type O blood type in the sample?

- Approach I (exact!):  $P(X=6)=?$  Where  $X \sim B(12, 0.44) \rightarrow$

$$P(X=6) = \binom{12}{6} (0.44)^6 (0.56)^6 = 0.2068 \text{ (SOCR)}$$

- Approach II (Approximate):  $X \sim B(n=12, p=0.44) \rightarrow$

$X$  (approx.)  $\sim N[\mu = n.p = 5.28; (np(1-p))^{1/2} = 1.7] \rightarrow P(X=6) \approx P(Z_1 \leq Z \leq Z_2)$ , where  $Z = (X - 5.28)/1.7$  and  $X_1=5.5, X_2=6.5$

So,  $P(X=6) \approx P(Z_1 \leq Z \leq Z_2) = 0.211$

Using Binomial Counts!

Slide 71

Stat 13, UCLA, Jon Dineen

### Example – 5.39 (in the textbook)

- Suppose proportion of blood type O is 0.44. If we take a random sample of 12 subjects and make a note of their blood types what is the probability that exactly 6 subjects have type O blood type in the sample?

- Approach I (exact!):  $P(X=6)=?$  Where  $X \sim B(12, 0.44) \rightarrow$

$$P(X=6) = \binom{12}{6} (0.44)^6 (0.56)^6 = 0.2068 \text{ (SOCR)}$$

- Approach III (Approximate):  $X \sim B(n=12, p=0.44) \rightarrow$

$p^* = X/n$  (approx.)  $\sim N[\mu = p = 0.44; (p(1-p)/n)^{1/2} = 0.1433] \rightarrow P(X=6) = P(p^* \approx 0.5) \approx P(p_1 \leq p^* \leq p_2)$ , where  $p_1 = 0.5 - 1/24$  and  $p_2 = 0.5 + 1/24$ . Standardize  $Z = (p - 0.44)/0.1433$  to get:

$P(X=6) \approx P(p_1 \leq p^* \leq p_2) = P(Z_1 \leq Z \leq Z_2) = 0.211$

Using Proportion!

Slide 72

Stat 13, UCLA, Jon Dineen

