

# AP Statistics CED 9.1 Daily Video 1 (Skill 1.A)

## Introducing Statistics – Do Those Points Align?

### What Will We Learn?

How can we determine if the slope of a sample regression line is consistent with random variation from a population regression model?

### How Faithful is Old Faithful?

The Old Faithful geyser is the most popular attraction in Yellowstone National Park. People travel from all over the world to see this geyser erupt. The National Park Service helps visitors plan their time in the park by when Old Faithful will erupt next. The Starnes family took its first trip to Yellowstone National Park in July 1995. They only had six hours in the park, but were able to see Old Faithful erupt. Mr. and Mrs. Starnes returned to the park in July 2019. They wondered if the model used to predict eruptions of Old Faithful was still the same as in 1995.

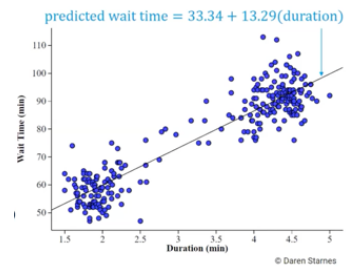
### Old Faithful: 1995

This scatterplot displays data on the duration (in minutes) and wait time until the next eruption (in minutes) for all 262 recorded eruptions of Old Faithful in July 1995.

There is a \_\_\_\_\_, \_\_\_\_\_ relationship between \_\_\_\_\_ and \_\_\_\_\_ in this population of Old Faithful eruptions.

Its equation is:

$$\text{predicted wait time} = 33.34 + 13.29 (\text{duration})$$



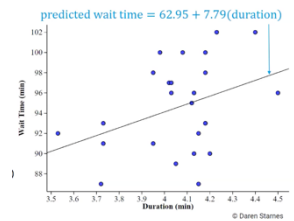
### Old Faithful: 2019

The scatterplot displays data on the duration (in minutes) and wait time until the next eruption (in minutes) for a random sample of 25 Old Faithful eruptions in July 2019.

There is a \_\_\_\_\_, \_\_\_\_\_ relationship between \_\_\_\_\_ and \_\_\_\_\_ in this population of Old Faithful eruptions.

Its equation is

$$\text{predicted wait time} = 62.95 + 7.79 (\text{duration})$$



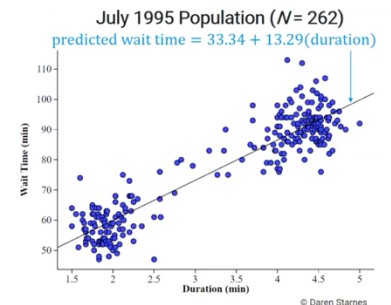
### Old Faithful: Then and Now

Is it believable that the population regression model from 1995 is still valid for predicting wait time from the duration of the previous Old Faithful eruption in 2019?

We need to \_\_\_\_\_ the \_\_\_\_\_ of getting a sample regression line with a \_\_\_\_\_ as least as unusual as 7.79 in a random sample of  $n = 25$  observations from the July 1995 population.

July 2019. Sample ( $n = 25$ ) yielded a:

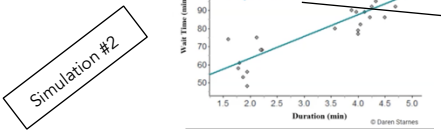
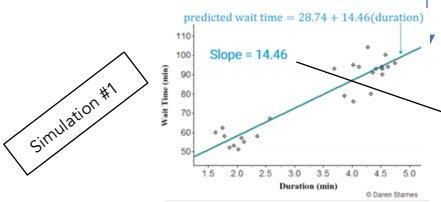
$$\text{predicted wait time} = 62.95 + 7.79(\text{duration})$$



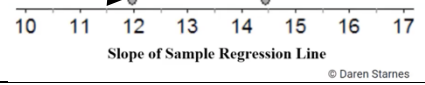
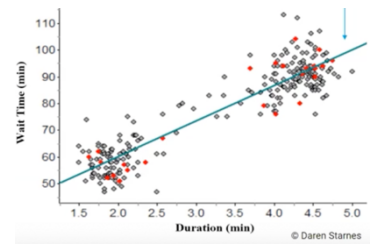
**Old Faithful: Simulation Trial 1**

We will simulate random samples of  $n =$  \_\_\_\_\_ points from the July 1995 population to see how likely it is we will get a slope of \_\_\_\_\_.

Simulated random sample of  $n=25$  points from the July 1995 population  
 predicted wait time =  $28.74 + 14.46(\text{duration})$



July 1995 Population ( $N = 262$ )  
 predicted wait time =  $33.34 + 13.29(\text{duration})$



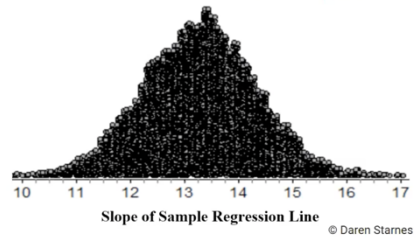
**Old Faithful: Simulation Results**

Is it believable that the population regression model from 1995 is still valid for predicting wait time from the duration of the previous Old Faithful eruption in 2019? After running many, many simulations we get the sampling distribution at the right.

July 2019 Sample ( $n = 25$ )  
 predicted wait time =  $62.95 + 7.79(\text{duration})$

We need to \_\_\_\_\_ the probability of getting a sample regression line with a \_\_\_\_\_ at least as unusual as \_\_\_\_\_ in a random sample of  $n = 25$  observations from the July, 1995 population.

July 1995 Population ( $N = 262$ )  
 predicted wait time =  $33.34 + 13.29(\text{duration})$   
 Simulated sampling distribution of the slope of the sample regression line based on samples of size  $n = 25$  from the July 1995 population



Looking at the sampling distribution, we see that we never got a slope of 7.79, so our estimated probability  $\approx$  \_\_\_\_\_

**Old Faithful: Simulation Results**

Is it believable that the population regression model from 1995 is still valid for predicting wait time from the duration of the previous Old Faithful eruption in 2019?

\_\_\_\_\_. There is an \_\_\_\_\_ 0 probability of obtaining a \_\_\_\_\_ regression line with a \_\_\_\_\_ as least as surprising (in \_\_\_\_\_ directions) as 7.79 is the population regression model from \_\_\_\_\_ is still \_\_\_\_\_.

**What Should We Take Away?**

- Take \_\_\_\_\_ random samples of size \_\_\_\_\_ from the population.
- Calculate the \_\_\_\_\_ of each sample regression line.
- Build the \_\_\_\_\_ distribution of the slope.
- See if the \_\_\_\_\_ value of the sample slope can be explained by \_\_\_\_\_ variation or not.

# AP Statistics CED 9.2 Daily Video 1 (Skill 1.D)

## Confidence Intervals for the Slope of a Regression Model

### What Will We Learn?

What conditions must the population regression model meet to obtain valid confidence intervals and significance tests for the slope?

How can we determine the shape, center and variability of the sampling distribution of the slope of a sample regression line?

### Is Old Faithful Still Faithful?

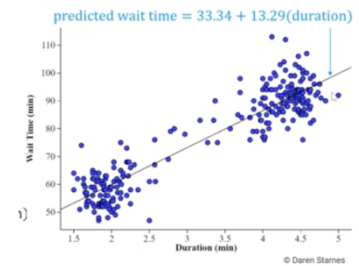
The Starnes family took its first trip to Yellowstone National Park in July 1995. They only had six hours in the park, but were able to see Old Faithful erupt. Mr. and Mrs. Starnes returned to the park in July 2019. They wondered if the model used to predict eruptions of Old Faithful was still the same as in 1995. Earlier, we used simulation to determine that the answer is “No.”

Can we construct a confidence interval for the slope of the population regression line in 2019? To answer that question, we need to understand the connection between the population regression model and the sampling distribution of the slope.

### Old Faithful: 1995

The scatterplot displays data on the duration (in minutes) and wait time until the next eruption (in minutes) for all 262 recorded eruptions of Old Faithful in July, 1995. We have added the population \_\_\_\_\_ line to scatterplot. Its equation is:

$$\text{predicted wait time} = 33.34 + 13.29(\text{duration})$$

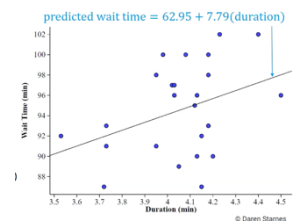


### Old Faithful: 2019

The scatterplot displays data on the duration (in minutes) and wait time until the next eruption (in minutes) for a random sample of 25 Old Faithful eruptions in July 2019.

We have added the \_\_\_\_\_ line to the scatterplot. Its equation is:

$$\text{predicted wait time} = 62.95 + 7.79(\text{duration})$$

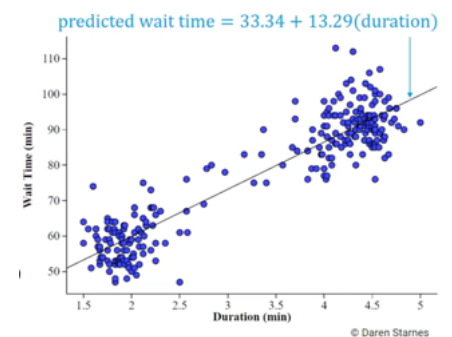


Can we construct a confidence interval for the slope of the population regression line in 2019?

### Simulated Sampling Distribution

Suppose we take \_\_\_\_\_ random sample of  $n = 25$  observations from the \_\_\_\_\_ of eruptions in July, 1995 and calculate the sample regression line  $\hat{y} = a + bx$  for each one.

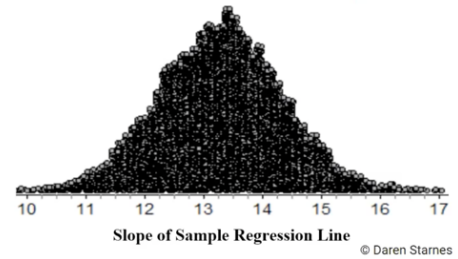
Can we determine the shape, center and variability of the sampling distribution of the slope  $b$  of the sample regression line from the population regression model?



**Simulated Sampling Distribution**

Suppose we take repeated random samples of  $n = 25$  observations from the population of eruptions in July 1995 and calculate the sample regression line  $\hat{y} = a + bx$  for each one. Can we determine the shape, center and variability of the sampling distribution of the slope  $b$  of the sample regression line from the population model?

July 1995 Population ( $N = 262$ )  
 predicted wait time =  $33.34 + 13.29(\text{duration})$   
 Simulated sampling distribution of the slope of the sample regression line based on samples of size  $n = 25$  from the July 1995 population



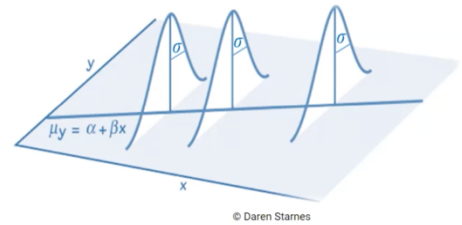
What are the shape, center and variability of the simulated sampling distribution of the slope of the sample regression line  $\hat{y} = a + bx$  ?

**Sampling Distribution of  $b$**  (Use the distribution above.)

Shape: \_\_\_\_\_ How are these characteristics related to the  
 Center: \_\_\_\_\_ population regression model?  
 Variability: \_\_\_\_\_

**Population Regression Model**

The population regression model is  $\mu_y = \alpha + \beta x$  where  $\mu_y$  is the mean value of the \_\_\_\_\_ variable  $y$  for a given value of the \_\_\_\_\_ variable  $x$ .



Confidence intervals and significance test for the slope  $\beta$  require that the population regression model meet these conditions:

- The \_\_\_\_\_ relationship between  $x$  and  $y$  is \_\_\_\_\_.
- The standard deviation of  $y$ ,  $\sigma_y$  does not \_\_\_\_\_ with  $x$ .
- For a particular value of  $x$ , the responses (\_\_\_\_\_) are approximately normally distributed.

**Population Regression Model**

Confidence intervals and significance test for the slope  $\beta$  require that the population regression model meet these conditions:

- The true relationship between  $x$  and  $y$  is linear. To check this:  
Scatterplot: The scatterplot shows a \_\_\_\_\_ relationship between duration and wait time.  
Residual Plot: The residual plot show \_\_\_\_\_ scatter about Residual = \_\_\_\_\_ line and no evidence of a \_\_\_\_\_ curved pattern.
- The standard deviation of  $y$ ,  $\sigma_y$ , does not vary with  $x$ .  
Scatterplot: The wait times \_\_\_\_\_ by a \_\_\_\_\_ amount for the different eruption durations in the data set.  
Residual Plot: The residuals are \_\_\_\_\_ in \_\_\_\_\_ for the eruption durations in the data set.
- For a particular value of  $x$ , the responses (\_\_\_\_\_) are approximately normally distributed.  
Dotplot: A dotplot of residuals is roughly \_\_\_\_\_, \_\_\_\_\_, and somewhat \_\_\_\_\_.



**Sampling Distribution of Slope**

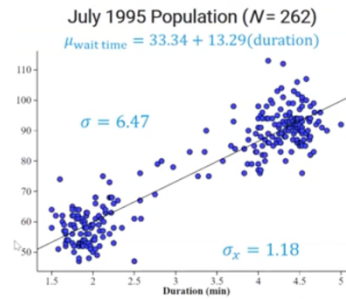
When a population regression model  $\mu_y = \alpha + \beta x$  meets the conditions, the sampling distribution of the slope  $b$  of the sample regression line has:

Shape: \_\_\_\_\_

Center:  $\mu_b = \beta =$  \_\_\_\_\_

Variability:  $\sigma_b = \frac{\sigma}{\sigma_x \sqrt{n}}$

where: \_\_\_\_\_ = Standard deviation of \_\_\_\_\_ for the population regression line and  
 \_\_\_\_\_ = Standard deviation of \_\_\_\_\_ in population and  
 \_\_\_\_\_ = sample size.



**Sampling Distribution of Slope**

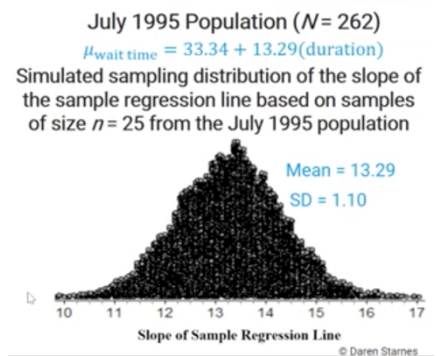
Suppose we take every possible random sample of \_\_\_\_\_ observations from the population and calculate the sample regression line  $\hat{y} = a + bx$  for each one.

**Sampling distribution of  $b$**

Shape: \_\_\_\_\_

Center:  $\mu_b = \beta =$  \_\_\_\_\_

Variability:  $\sigma_b = \frac{\sigma}{\sigma_x \sqrt{n}} =$  \_\_\_\_\_



When sampling without replacement check: (if  $n \leq 10\%N$ ) \_\_\_\_\_

**Inference for Slope**

Can we construct a confidence interval for the slope of the population regression line in 2019?

Related questions:

How do we check the conditions about the \_\_\_\_\_ regression model using only the \_\_\_\_\_ data?

How can we \_\_\_\_\_ the standard deviation of the \_\_\_\_\_ distribution of the sample slope,  $\sigma_b = \frac{\sigma}{\sigma_x \sqrt{n}}$ , when we do not know the values of  $\sigma$  or  $\sigma_x$ ?

**What Should We Take Away?**

What conditions must the population regression model meet to obtain valid confidence intervals and significance tests for the slope?

- The \_\_\_\_\_ relationship between  $x$  and  $y$  is \_\_\_\_\_.
- The standard deviation of  $y$ ,  $\sigma_y$ , does \_\_\_\_\_ with  $x$ .
- For a particular value of  $x$ , the responses \_\_\_\_\_ are approximately normally distributed.

How can we determine the shape, center and variability of the sampling distribution of the slope of a sample regression line?

Shape: Approximately normal;

Center:  $\mu_b = \beta$

Variability:  $\sigma_b = \frac{\sigma}{\sigma_x \sqrt{n}}$

## AP Statistics CED 9.2 Daily Video 2 (Skill 4.C)

### Confidence Intervals for the Slope of a Regression Model

#### What Will We Learn?

How do we identify an appropriate confidence interval procedure for the slope of a population regression line?

How do we verify the conditions for calculating a confidence interval for the slope of a population regression line?

#### Predicting Old Faithful in 2019

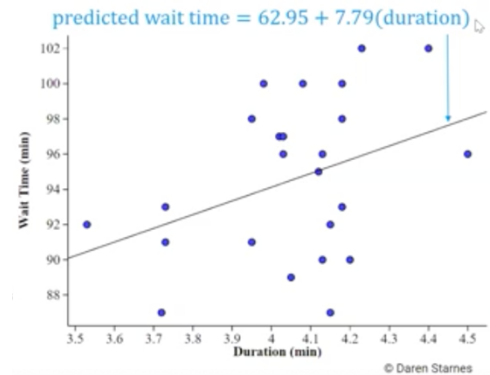
The Starnes family took its first trip to Yellowstone National Park in July 1995. They only had six hours in the park, but were able to see Old Faithful erupt. Mr. and Mrs. Starnes returned to the park in July 2019. They wondered if the model used to predict eruptions of Old Faithful was still the same as in 1995. Earlier, we used simulation to determine that the answer is “No.”

Can we construct a confidence interval for the slope of the population regression line in 2019?

#### Old Faithful: 2019

The scatterplot displays data on the duration (in minutes) and wait time until the next eruption (in minutes) for a random sample of 25 Old Faithful eruptions in July 2019, along with the sample regression line. Computer output from a least-squares regression analysis is shown below:

Predictor	Coef	SE Coef	T-Value	P-Value
Constant	62.95	16.4	3.85	0.001
Duration(min)	7.79	4.03	1.94	0.065
S=4.20970		R-sq=14.01%		R-sq(adj)=10.27%



Can we construct a confidence interval for the slope of the population regression line in 2019?

#### Population Regression Model

The population regression model is  $\mu_y = \alpha + \beta x$  where  $\mu_y$  is the mean value of the \_\_\_\_\_ variable  $y$  for a given value of the \_\_\_\_\_ variable  $x$ . Confidence intervals and significance tests for the slope  $\beta$  require that the population regression model meet three conditions:

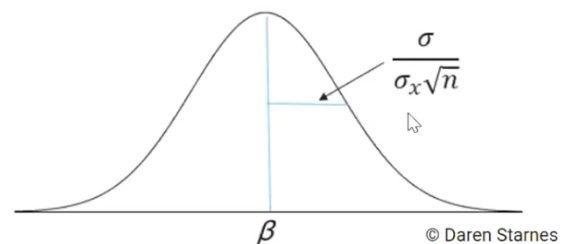
- The \_\_\_\_\_ relationship between  $x$  and  $y$  is \_\_\_\_\_.
- The standard deviation of  $y$ ,  $\sigma_y$  does not \_\_\_\_\_ with  $x$ .
- For a particular value of  $x$ , the responses (\_\_\_\_\_) are approximately normally distributed.

#### Sampling Distribution of Slope

When the population regression model  $\mu_y = \alpha + \beta x$  meets the conditions, the sampling distribution of the slope  $b$  of the sample regression line has:

Shape: \_\_\_\_\_;  
 Center: \_\_\_\_\_; and  
 Variability \_\_\_\_\_ (\_\_\_\_\_)

Where: \_\_\_\_\_ = standard deviation of residuals for the population regression line  
 \_\_\_\_\_ = standard deviation of  $x$ -values in population and  
 \_\_\_\_\_ = sample size.



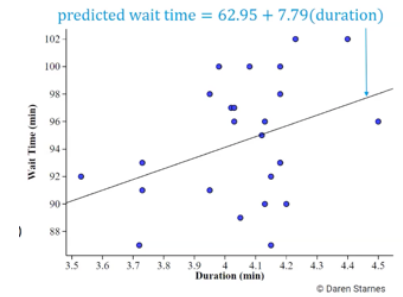
### Inference for Slope

Can we construct a confidence interval for the slope of the population regression line in 2019?

#### Related questions:

How do we check the conditions about the \_\_\_\_\_ regression model using only the \_\_\_\_\_ data?

How can we \_\_\_\_\_ the standard deviation of the \_\_\_\_\_ distribution of the sample slope,  $\sigma_b = \frac{\sigma}{\sigma_x \sqrt{n}}$ , when we do not know the values of  $\sigma$  or  $\sigma_x$ ?



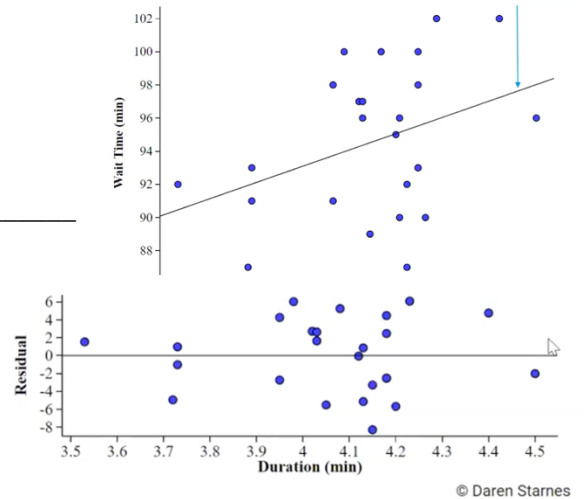
### Checking Conditions

To calculate a confidence interval for the slope  $\beta$  of a population regression line we must check the following:

- **The true relationship between  $x$  and  $y$  is linear.**

Scatterplot: The scatterplot shows a \_\_\_\_\_ relationship between duration and wait time.

Residual Plot: The residual plot shows \_\_\_\_\_ scatter about Residual = \_\_\_\_\_ line and no evidence of a \_\_\_\_\_ curved pattern.

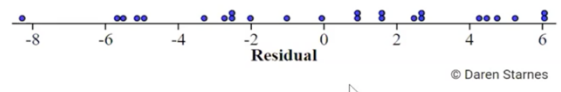


- **The standard deviation of  $y$ ,  $\sigma_y$ , does not vary with  $x$ .**

Scatterplot: The wait times \_\_\_\_\_ by a \_\_\_\_\_ amount for the different eruption durations in the data set.

Residual Plot: The residuals are \_\_\_\_\_ in \_\_\_\_\_ for the eruption durations in the data set.

- **For a particular value of  $x$ , the responses (\_\_\_\_\_ ) are approximately normally distributed.**



Dotplot: A dotplot of residuals show no \_\_\_\_\_.

If the observed distribution is \_\_\_\_\_ or shows other \_\_\_\_\_ from normality,  $n$  should be \_\_\_\_\_ or we cannot proceed.

- **There is independence in data collection.** Data are collected using a \_\_\_\_\_ sample or a \_\_\_\_\_ experiment.

In this case, the data came from a \_\_\_\_\_ sample of \_\_\_\_\_ Old Faithful eruptions in \_\_\_\_\_.

When sampling without replacement, check the  $n \leq 10\%N$ . In this case, \_\_\_\_\_  $\leq$  \_\_\_\_\_

All of the conditions are \_\_\_\_\_

### Inference for Slope

Back to our question, "How can we estimate the standard deviation of the sampling distribution of the sample slope  $\sigma_b = \frac{\sigma}{\sigma_x \sqrt{n}}$ , when we do not know the values of  $\sigma$  or  $\sigma_x$ ?"

**Estimating  $\sigma$**

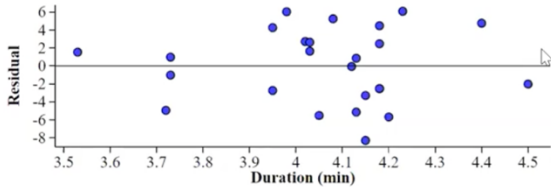
We can use the sample regression \_\_\_\_\_.

The \_\_\_\_\_ standard deviation of the x-values is  $s_x =$  \_\_\_\_\_. We can use  $s_x$  to estimate  $\sigma_x$ .

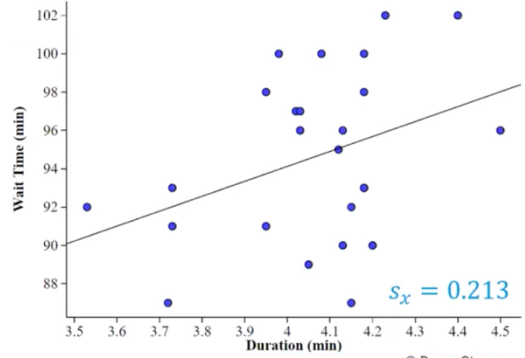
The standard deviation of the residuals for the sample regression line:

$$s = \sqrt{\frac{\sum (y_i - \hat{y})^2}{n - 2}} = \sqrt{\frac{\sum \text{residual}^2}{n - 2}}$$

$s$  estimates the size of the typical prediction error



© Daren Starnes



$s_x = 0.213$

© Daren Starnes

Predictor	Coef	SE Coef	T-Value	P-Value
Constant	62.95	16.4	3.85	0.001
Duration (min)	7.79	4.03	1.94	0.065
S=4.20970		R-sq=14.01%	R-sq(adj)=10.27%	

From the output, we can see that  $s =$  \_\_\_\_\_. We use this value to estimate  $\sigma$ .

**Confidence Interval for Slope**

Can we construct a confidence interval for the slope of the population regression line in 2019? YES!

The appropriate confidence interval is a \_\_\_\_\_.

**What Should We Take Away?**

How do we identify an appropriate confidence interval procedure for the slope of a population regression line?

The appropriate confidence interval is a \_\_\_\_\_ because we are \_\_\_\_\_  $\sigma$  with  $s$ .

How do we verify the conditions for calculating a confidence interval for the slope of a population regression line?

- \* The true relationship between \_\_\_\_\_ is linear. With \_\_\_\_\_ and \_\_\_\_\_
- \* The standard deviation of  $y$ ,  $\sigma_y$ , does vary with  $x$ . With \_\_\_\_\_ and \_\_\_\_\_
- \* For a particular value of  $x$ , the  $y$ -values are approximately normally distributed. With \_\_\_\_\_ of the \_\_\_\_\_. Need \_\_\_\_\_ if obvious \_\_\_\_\_ or \_\_\_\_\_.
- \* There is \_\_\_\_\_ in data collection  
 Data are collected using a \_\_\_\_\_ sample or a \_\_\_\_\_ experiment.  
 When sampling \_\_\_\_\_ replacement, check the \_\_\_\_\_.

# AP Statistics CED 9.2 Daily Video 3 (Skill 3.D)

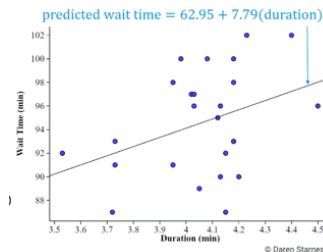
## Confidence Intervals for the Slope of a Regression Model

### What Will We Learn?

How do we determine the margin of error when estimating the slope of a population regression line?  
 How do we calculate a confidence interval for the slope of a population regression line?

### Predicting Old Faithful: 2019

The scatterplot displays data on the duration (in minutes) and wait time until the next eruption (in minutes) for a random sample of 25 Old Faithful eruptions in July 2019, along with the sample regression line.



Computer output from a least-square regression analysis is shown below.

Predictor	Coef	SE Coef	T-Value	P-Value
Constant	62.95	16.4	3.85	0.001
Duration(min)	7.79	4.03	1.94	0.065
S=4.20970		R-sq=14.01%		R-sq(adj)=10.27%

Calculate and interpret a 95% confidence interval for the slope of the population regression line in 2019.

### Calculating the Margin of Error

In AP Statistics, confidence intervals have the form:

$$CI = \text{_____} \pm \text{_____}$$

The margin of error describes how \_\_\_\_\_ a value of a \_\_\_\_\_ statistic is \_\_\_\_\_ to vary from the value of the \_\_\_\_\_ population \_\_\_\_\_.

The margin of error is determined by \_\_\_\_\_ factors:

- How much the statistic \_\_\_\_\_ varies from the \_\_\_\_\_.
- How \_\_\_\_\_ we want to be in our \_\_\_\_\_.

$$\text{margin of error} = (\text{_____})(\text{_____})$$

### Standard Error of Slope

The standard error of a statistic is an \_\_\_\_\_ of the \_\_\_\_\_ of the sampling distribution of the \_\_\_\_\_.

From topic 9.2, the standard deviation of the sampling distribution of  $b$  is:

$$\sigma_b = \frac{\sigma}{\sigma_x \sqrt{n}}$$

Because we don't know the value of  $\sigma$  or  $\sigma_x$ , we estimate them using the standard deviation of \_\_\_\_\_ for the sample regression line,  $s$ , and the standard deviation of the x-values in the sample,  $s_x$ , to get the \_\_\_\_\_:

$$SE_b = \frac{s}{s_x \sqrt{n - 1}}$$

### Calculating the Margin of Error

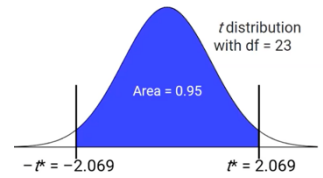
The critical value is a \_\_\_\_\_ that makes the margin of error large enough to give a specific amount of confidence that the \_\_\_\_\_ contains the value of the \_\_\_\_\_.

From confidence intervals for the \_\_\_\_\_ of a population regression line, the \_\_\_\_\_ represent the \_\_\_\_\_ encompassing the middle \_\_\_\_\_ of the \_\_\_\_\_ with degrees of freedom \_\_\_\_\_, where C% is the \_\_\_\_\_.

**Calculating the Margin of Error**

\_\_\_\_\_ = (critical value)(standard error of statistic)

In the Old Faithful example, we are asked to construct a 95% confidence interval based on a random sample of  $n = 25$  eruptions. To find the critical value \_\_\_\_\_ for a 95% confidence interval, find the \_\_\_\_\_ encompassing the middle \_\_\_\_\_ of the \_\_\_\_\_ with \_\_\_\_\_. This value can be found using Table B or using technology using InvT. Either way the critical value  $t^* =$  \_\_\_\_\_



**Calculating the Confidence Interval**

CI = \_\_\_\_\_  $\pm$  \_\_\_\_\_

CI = \_\_\_\_\_  $\pm$  ( \_\_\_\_\_ )( \_\_\_\_\_ )

CI =  $b \pm t^* SE_b = b \pm t^* \frac{s}{s_x \sqrt{n-1}}$

(Use this space to calculate the CI)

For our example about Old Faithful eruptions,  $b =$  \_\_\_\_\_ and  $s =$  \_\_\_\_\_. Unfortunately,  $s_x$  is not shown in the computer regression output. From the previous video,  $s_x =$  \_\_\_\_\_.

Predictor	Coef	SE Coef	T-Value	P-Value
Constant	62.95	16.4	3.85	0.001
Duration (min)	7.79	4.03	1.94	0.065
S=4.20970		R-sq=14.01%	R-sq(adj)=10.27%	

↑ SD of residuals

↑ standard error of slope

**Calculating the Confidence Interval**

All of the components needed to calculate a confidence interval can be found on the AP Statistics formula sheet. Make sure you can locate the information as you watch the video.

**Factors that Affect Interval Width**

Recall that confidence intervals in AP Statistics have the following structure:

CI = \_\_\_\_\_  $\pm$  \_\_\_\_\_

The width of a confidence interval is \_\_\_\_\_ the \_\_\_\_\_.

For a confidence interval about the slope of a population \_\_\_\_\_,

\_\_\_\_\_ = \_\_\_\_\_

We generally prefer \_\_\_\_\_ confidence intervals ( \_\_\_\_\_ ), so we want the margin of error to be \_\_\_\_\_. There are two common ways to decrease margin of error.

(1) \_\_\_\_\_ the \_\_\_\_\_ and (2) \_\_\_\_\_ the \_\_\_\_\_.

**Calculating a Confidence Interval**

Raoul performed an experiment using 16 windup rubber band single-propellor airplanes. He wound up the propeller a different number of times and recorded the amount of time (in seconds) that the airplane flew for each number of rotations that the propeller was wound. A regression analysis was performed and the partial computer output is given below. Assuming that the conditions for inference are satisfied and calculate a 95% CI.

Predictor	Coef	SE Coef	T	P
Constant	0.9241	0.6413	1.44	0.172
Rotation	0.04625	0.01565	2.96	0.010
S = 0.5426		R-Sq = 38.4%	R-Sq(adj) = 34.0%	



Assuming that the conditions for inference are satisfied, which of the following is a 95 percent confidence interval for the slope of the regression line that relates to the number of rotations the rubber band is wound and plane's flight time?

(A)  $0.04625 \pm (2.145)(0.01565)$

(B)  $0.9241 \pm (2.145)(0.6413)$

(C)  $0.04625 \pm (2.131)(0.01565)$

(D)  $0.04625 \pm (2.131)\left(\frac{0.5426}{\sqrt{16}}\right)$

(E)  $0.9241 \pm (2.131)(0.6413)$

$$CI = b \pm t^* \frac{s}{s_x \sqrt{n-1}} = b \pm t^* SE_b$$

### What Should We Take Away?

How do we determine the margin of error when estimating the slope of a population regression line?

$$\text{margin of error} = (\text{critical value})(\text{standard error of statistic})$$

$$\text{margin of error} = \text{critical value} \times \text{standard error of statistic}$$

How do we calculate a confidence interval for the slope of a population regression line?

$$CI = \text{point estimate} \pm \text{margin of error}$$

$$CI = b \pm t^* SE_b = b \pm t^* \frac{s}{s_x \sqrt{n-1}}$$

## AP Statistics CED 9.3 Daily Video 1 (Skill 4.B)

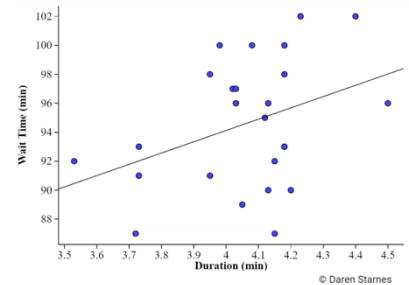
### What Will We Learn?

How do we interpret a confidence interval for the slope of a population regression line?  
How do we justify a claim based on a confidence interval for slope?

### Old Faithful: 2019

The scatterplot displays data on the duration (in minutes) and wait time until the next eruption (in minutes) for a random sample of 25 Old Faithful eruptions in July 2019, along with the sample regression line. Computer output from the least-squares regression analysis is shown below:

Predictor	Coef	SE Coef	T-Value	P-Value
Constant	62.95	16.4	3.85	0.001
Duration (min)	7.79	4.03	1.94	0.065
S=4.20970		R-sq=14.01%		R-sq (adj)=10.27%



Calculate and interpret a 95% confidence interval for the slope of the population regression line in 2019.

### Interpreting the Confidence Interval

In general, here is how to interpret a confidence interval for the slope of a regression model:

"We are \_\_\_\_\_ confident that the \_\_\_\_\_ from \_\_\_\_\_ to \_\_\_\_\_ \_\_\_\_\_ the slope of the population regression line [\_\_\_\_\_]."

From Topic 9.2 Video 3, the 95% confidence interval is -0.55 to 16.13:

"We are \_\_\_\_\_ confident that the \_\_\_\_\_ from \_\_\_\_\_ to \_\_\_\_\_ \_\_\_\_\_ the slope of the \_\_\_\_\_ regression line for predicting \_\_\_\_\_ until the next eruption (in minutes) from the \_\_\_\_\_ of the previous eruption (in minutes) for \_\_\_\_\_ Old Faithful geyser eruptions in July 2019."

### Interpreting Confidence Level

In \_\_\_\_\_ random sampling with the \_\_\_\_\_ sample size, approximately C% of "C%" confidence intervals created \_\_\_\_\_ the slope of the \_\_\_\_\_ regression line. If we take \_\_\_\_\_ random samples of size \_\_\_\_\_ from the population of Old Faithful eruptions in \_\_\_\_\_, and use each sample to \_\_\_\_\_ a 95% confidence interval for the slope of the \_\_\_\_\_ regression line for predicting \_\_\_\_\_ until the next eruption (in minutes) from the \_\_\_\_\_ of the previous eruption (in minutes), about \_\_\_\_\_ of those \_\_\_\_\_ would \_\_\_\_\_ the \_\_\_\_\_ slope.

### Justifying a Claim

Does the confidence interval (-0.55 to 16.13) provide convincing evidence that wait time until the next eruption of Old Faithful is linearly related to the duration of the previous eruption in July 2019?

Note that  $\beta = 0$  would indicate a line with a \_\_\_\_\_ is the model for \_\_\_\_\_ wait time from eruption duration. In other words, the regression model would predict the \_\_\_\_\_ wait time until the next Old Faithful eruption no matter the duration of the previous eruption.

Because the confidence interval (\_\_\_\_\_) contains \_\_\_\_\_ as a \_\_\_\_\_ value of the slope of the \_\_\_\_\_ regression line, there is \_\_\_\_\_ evidence that \_\_\_\_\_ until the next eruption of Old Faithful is \_\_\_\_\_ related to the \_\_\_\_\_ of the previous eruption in July 2019.

**Interpreting a Confidence Interval**

Raoul performed an experiment using 16 windup rubber band single-propellor airplanes. He wound up the propeller a different number of times and recorded the amount of time (in seconds) that the airplane flew for each number of rotations that the propeller was wound. A regression analysis was performed, and the conditions for inference were verified. A 95% confidence interval for the slope of the regression line that relates the number of rotations the rubber band is wound and the plane's flight time is given by (0.013, 0.080). Which of the following provides a correct interpretation of the confidence interval?

- (A) There is a 0.95 probability that the slope of the population regression line that relates the number of rotations the rubber band is wound and the plane's flight time is between 0.013 and 0.080.
- (B) If the data collection process were repeated many times, about 95% of the resulting sample regression lines would have slopes between 0.013 and 0.080.
- (C) If the data collection process were repeated many times, about 95% of the resulting confidence intervals would contain the slope of the sample regression line.
- (D) We are 95% confident that the slope of the sample regression line that relates the number of rotations the rubber band is wound and the plane's flight time is between 0.013 and 0.080.
- (E) We are 95% confident that the slope of the population regression line that relates the number of rotations the rubber band is wound and the plane's flight time is between 0.013 and 0.080.

**Justifying a Claim with a CI**

- (A) There is not convincing evidence of a linear relationship between the number of rotations of the rubber band and the flight times of these 16 windup airplanes because 0 is not included in the interval.
- (B) There is not convincing evidence of a linear relationship between the number of rotations of the rubber band and the flight times of windup airplanes like these because 0 is not included in the interval.
- (C) There is convincing evidence of a positive linear relationship between the number of rotations of the rubber band and the flight times of these 16 windup airplanes because all values in the interval are positive.
- (D) There is convincing evidence of a positive linear relationship between the number of rotations of the rubber band and the flight times of windup airplanes like these because all values in the interval are positive.
- (E) There is convincing evidence that there is not a linear relationship between the number of rotations of the rubber band and the flight times of windup airplanes like these because 0 is not included in the interval.

**What Should We Take Away?**

How do we interpret a confidence interval for the slope of a population regression line?

"We are \_\_\_\_\_ confident that the \_\_\_\_\_ from \_\_\_\_ to \_\_\_\_ \_\_\_\_\_ the slope of the population regression line [\_\_\_\_\_]."

How do we justify a claim based on a confidence interval for slope?

- If all the values in the confidence interval are \_\_\_\_\_ with the claim, there \_\_\_\_\_ convincing evidence for the claim.
- If one or more of the values in the confidence interval are \_\_\_\_\_ with the claim, there is \_\_\_\_\_ convincing evidence for the claim.

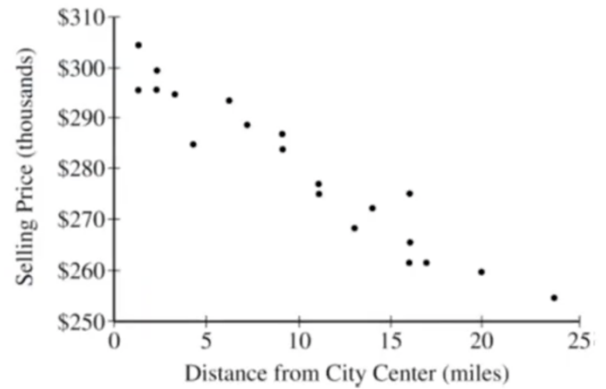
## AP Statistics CED 9.3 Daily Video 2 (Skill 3.E)

### What Will We Learn?

How do we construct and interpret a confidence interval for the slope of a population regression line?

#### 2019 International Exam #2

A real estate agent working in a large city believes that, for three-bedroom houses, the selling price of the house decreases by approximately \$2,000 for every mile increase in the distance of the house from the city center. To investigate the belief, the agent obtained a random sample of 20 three-bedroom houses that sold in the last year. The selling price, in thousands of dollars and the distance from the city center, in miles, for each of the 20 houses are shown in the scatterplot. The table shows computer output from a regression analysis of the data.



Predictor	Coef	SE Coef	T	P
Constant	301.7	1.80	167.17	0.000
Distance	-2.158	0.149	-14.45	0.000
S = 4.4336		R-sq = 92.1%		

(a) Assume all conditions for inference are met.

Construct and interpret a 95 percent confidence interval for the slope of the least-squares regression line.

(b) Does the confidence interval contradict the agent's belief about the relationship between selling price and distance from the city center? Justify your answer.

#### 2019 International Exam #2(a) Calculate Interval

(a) Assume all conditions for inference are met. Construct and interpret a 95 percent confidence interval for the slope of the least-squares regression line.

We will construct a \_\_\_\_\_ confidence interval for  $\beta =$  \_\_\_\_\_ of the population regression line for \_\_\_\_\_ selling price (in thousands) from distance from the city center (in miles) for all three-bedroom houses near this city.

We will use a \_\_\_\_\_ and fortunately conditions are \_\_\_\_\_.  
(Use the output above to create the t-interval.)

$b =$  \_\_\_\_\_ (Find on Output)       $df =$  \_\_\_\_\_       $t^* =$  \_\_\_\_\_ (use technology InvT)

Formula:  $CI = b \pm t^*SE_b$   
 = \_\_\_\_\_  
 = \_\_\_\_\_  
 = \_\_\_\_\_

#### 2019 International Exam #2, Part (a) Interpret Interval      Given: 95% CI = -2.471 to -1.845

We are \_\_\_\_\_ that the slope of the \_\_\_\_\_ regression line is between \_\_\_\_\_ and \_\_\_\_\_ thousands of \_\_\_\_\_. This implies that for \_\_\_\_\_ additional mile that a three-bedroom house is away from the \_\_\_\_\_, the selling price of the house is expected to \_\_\_\_\_ between \_\_\_\_\_ and \_\_\_\_\_.

**2019 International Exam #2**

A real estate agent working in a large city believes that, for three-bedroom houses, the selling price of the house decreases by approximately \$2,000 for every mile increase in the distance of the house from the city center. To investigate the belief, the agent obtained a random sample of 20 three-bedroom houses that sold in the last year. The selling price, in thousands of dollars and the distance from the city center, in miles, for each of the 20 houses are shown in the scatterplot. The table shows computer output from a regression analysis of the data.

(b) Does the confidence interval contradict the agent's belief about the relationship between selling price and distance from the city center? Justify your answer.

Given: 95% CI = -2.471 to -1.845, Interpret the interval

Because the confidence interval contains \_\_\_\_\_, corresponding to a \_\_\_\_\_ decrease, it is a \_\_\_\_\_ value for the slope of the regression line. Consequently, the data \_\_\_\_\_ contradict the agent's belief that the selling prices of three-bedroom houses \_\_\_\_\_ about \_\_\_\_\_ for every \_\_\_\_\_ increase in distance of the house from the city center.

**What Should We Take Away?**

How do we construct and interpret a confidence interval for the slope of a population regression line?

**Make sure to:**

- Define the \_\_\_\_\_ you are trying to estimate.
- Identify the \_\_\_\_\_ you are using.
- Verify that the \_\_\_\_\_ for the procedure are \_\_\_\_\_.
- \_\_\_\_\_ the confidence interval.
- \_\_\_\_\_ the interval in \_\_\_\_\_.

# AP Statistics CED 9.4 Daily Video 1 (Skill 1.F)

## Setting Up a Test for the Slope of a Regression Model

### What Will We Learn?

How do you state a null hypothesis in a test about the slope of a population regression line?  
 How do you state an alternative hypothesis in a test about the slope of a population regression line?

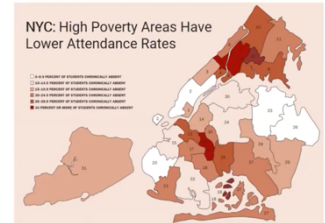
### Schools and "Equal Opportunity"

Can education systems equalize opportunities for lower income students?

### The Income Attendance Gap

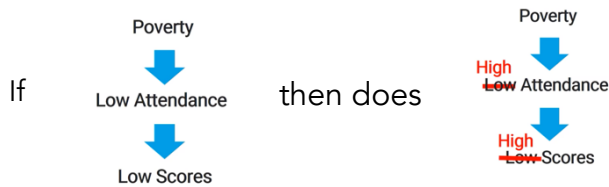
Nationally, higher income areas tend to have fewer chronically absent students. Possible reasons:

- Transportation access
- Work to support family



Graphic from Neuser et al., "A Better Picture of Poverty," Center for New York City Affairs, Nov. 2014. [https://www.attendanceworks.org/wp-content/uploads/2017/06/BetterPicturePoverty\\_A9\\_F2014\\_201.pdf](https://www.attendanceworks.org/wp-content/uploads/2017/06/BetterPicturePoverty_A9_F2014_201.pdf)

### Is Attendance the Solution?

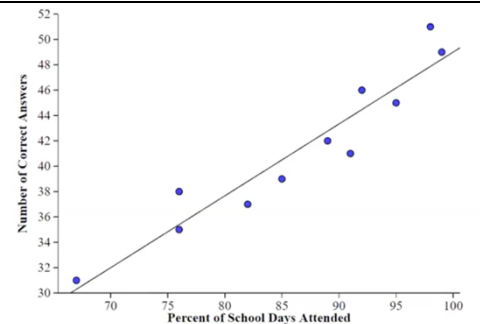


Some school systems have targeted attendance as the key to raising test scores for lower income students.

### Let's Look at the Data

Researchers collected data on the percent of school days attended and the number of questions answered correctly on the state's end-of-course test for a random sample of 11 Texas Algebra 1 students. Here are the data, along with a scatterplot and computer output from a least-squares regression analysis.

Percent attendance	95	89	67	98	99	76	92	91	76	85	82
Questions correct	45	42	31	51	49	38	46	41	35	39	37



Predictor	Coef	SE Coef	T	P
Constant	-7.69	5.37	-1.43	0.186
Attendance	0.57	0.062	9.18	0.000
S	1.99	R-Sq = 90.3%	R-Sq (adj) = 89.3%	

Do the data give convincing evidence at the  $\alpha = 0.01$  significance level of a positive, linear relationship between state test score and percent attendance for Texas Algebra 1 students?

### Null Hypothesis

In a statistical test, the \_\_\_\_\_ hypothesis is often a claim of "\_\_\_\_\_", "\_\_\_\_\_", or "\_\_\_\_\_". In the attendance and test scores example, the null hypothesis is that there is \_\_\_\_\_ between the percent of schools days attended and the number of questions answered correctly on the state test by Texas Algebra I students. In other words, the \_\_\_\_\_ regression line would have a \_\_\_\_\_. In symbols:

$H_0: \beta = \underline{\hspace{1cm}}$ ; where  $\beta = \underline{\hspace{1cm}}$  of the \_\_\_\_\_ regression line for \_\_\_\_\_ number of questions answered correctly on the state test from the percent of schools day attended for Texas Algebra 1 students.

Until we have \_\_\_\_\_ otherwise, we assume the \_\_\_\_\_ hypothesis is correct.



**Alternative Hypothesis**

In a statistical test, the \_\_\_\_\_ hypothesis is the \_\_\_\_\_ that we hope to support with \_\_\_\_\_ from the data collected. In the attendance and test scores example, the researchers wanted to know if there is a \_\_\_\_\_, \_\_\_\_\_ relationship between the \_\_\_\_\_ of school days attended by Texas students and the \_\_\_\_\_ of questions answered correctly on the state Algebra 1 test. So the alternative hypothesis is that the population regression line would have a \_\_\_\_\_. In symbols:

$H_0: \beta > \text{_____}$ ; where  $\beta = \text{_____}$  of the population regression line for predicting number of questions answered correctly on the state test from the percent of school days attended for Texas Algebra 1 students.

**Stating Hypotheses: Summary**

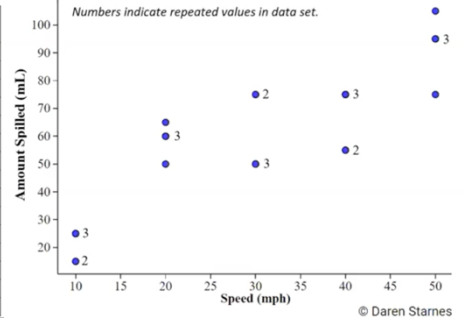
For hypotheses about the \_\_\_\_\_ of a population regression line:

- o The null is a statement of \_\_\_\_\_, typically \_\_\_\_\_.
- o The alternative always contains a strict \_\_\_\_\_, typically \_\_\_\_\_.  
When the inequality is \_\_\_\_\_, the alternative is called \_\_\_\_\_.  
When the inequality is \_\_\_\_\_, the alternative is called \_\_\_\_\_.  
The choice of alternative is determined by the \_\_\_\_\_ of interest and should be stated \_\_\_\_\_ data collection begins.
- o Never refer to \_\_\_\_\_ (such as *b*) in the hypotheses!
- o Remember to \_\_\_\_\_ the \_\_\_\_\_.

**Don't Spill my Drink!**

Two AP Statistics students wondered if there is a linear relationship between speed and amount of drink spilled when driving on bumpy dirt roads. To find out, they filled a cup with 275 mL of water and placed it in the car's cup holder, then drove down a bumpy road at a specified speed and recorded how much water spilled out. The 25 trials were randomly assigned to be 5 trials each at 10, 20, 30, 40, or 50 miles per hour (mph).

Speed (mph)	Spilled (mL)	Speed (mph)	Spilled (mL)
10	25	40	75
10	25	40	75
10	25	40	75
10	15	40	55
10	15	40	55
20	60	50	95
20	60	50	95
20	65	50	95
20	60	50	75
20	50	50	105
30	50		
30	75		
30	50		
30	75		
30	50		



Term	Coef	SE Coef
Constant	14.40	5.95
Speed (mph)	1.520	0.179
S=12.6800		R-sq = 75.75%

Here are the data. (Note that several data points have the same values for speed and amount spilled.) A scatterplot of the data, along with computer output from a least-squares regression analysis are shown. Do the data provide convincing evidence at the 0.05 significance level of a linear relationship between the car's speed on a bumpy dirt road and the amount of drink spilled? Assume the conditions for inference are met.

**Stating the Hypotheses**

$H_0: \text{_____}$  and  $H_a: \text{_____}$ ; where \_\_\_\_\_

**What Should We Take Away?**

How do you state a null hypothesis in a test about the slope of a population regression line?

$H_0: \beta = \beta_0$  For a null hypothesis of \_\_\_\_\_ linear relationship, \_\_\_\_\_

**Note: remember to clearly define the \_\_\_\_\_  $\beta$ .**

How do you state an alternative hypothesis in a test about the slope of a population regression line?

$H_a: \beta > \beta_0$  For a test of a positive linear relationship, \_\_\_\_\_. For the test of a negative

$H_a: \beta < \beta_0$  relationship \_\_\_\_\_.

$H_a: \beta \neq \beta_0$  relationship \_\_\_\_\_.

# AP Statistics CED 9.4 Daily Video 2 (Skill 4.C)

## Setting Up a Test for the Slope of a Regression Model

### What Will We Learn?

How do we identify an appropriate significance test procedure for the slope of a population regression line?

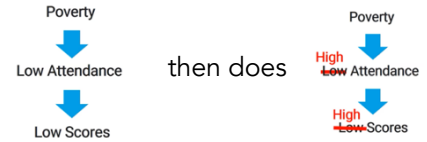
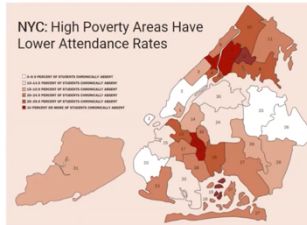
How do we verify the conditions for performing a test about the slope of a population regression line?

### Is Attendance the Solution?

Nationally, higher income areas tend to have fewer chronically absent students.

Possible reasons:

- Transportation access
- Work to support family

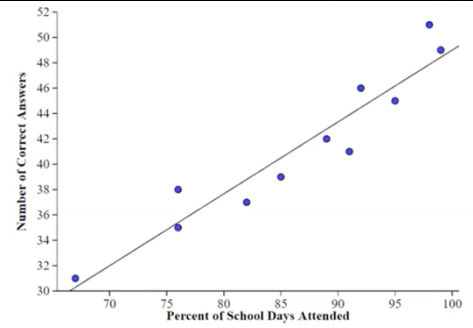


Some school systems have targeted attendance as the key to raising test scores for lower income students.

### Let's Look at the Data

Researchers collected data on the percent of school days attended and the number of questions answered correctly on the state's end-of-course test for a random sample of 11 Texas Algebra 1 students. Here are the data, along with a scatterplot and computer output from a least-squares regression analysis.

Percent attendance	95	89	67	98	99	76	92	91	76	85	82
Questions correct	45	42	31	51	49	38	46	41	35	39	37



Predictor	Coef	SE Coef	T	P
Constant	-7.69	5.37	-1.43	0.186
Attendance	0.57	0.062	9.18	0.000

S = 1.99    R-Sq = 90.3%    R-Sq (adj) = 89.3%

Do the data give convincing evidence at the  $\alpha = 0.01$  significance level of a positive, linear relationship between state test score and percent attendance for Texas Algebra 1 students?

### Test Scores and Attendance - Hypotheses

In a previous video, we stated the hypotheses:  $H_0: \beta = 0$  and  $H_a: \beta > 0$ ; where  $\beta$  = the slope of the population regression line for predicting number of questions answered correctly on the state test from the percent of schools day attended for Texas Algebra 1 students.

### Identifying the Procedure

You have learned many different significance test procedures this year. Some involved one sample and others involve two samples. Some involved inference about categorical data: \_\_\_\_\_ or \_\_\_\_\_. Others involved inference about quantitative data: \_\_\_\_\_ or \_\_\_\_\_.

When the goal is to test a claim about the \_\_\_\_\_ of a \_\_\_\_\_ regression line, we use a \_\_\_\_\_.

### Checking Conditions

To perform a significance test about the slope of  $\beta$  of a population regression line, we must check the following:

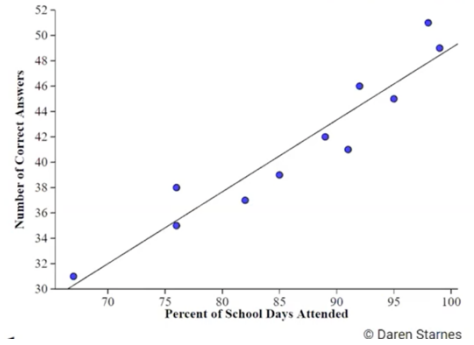
- The true relationship between \_\_\_\_\_ is linear.
- The standard deviation of \_\_\_\_\_ does not vary with \_\_\_\_\_/
- For a particular value of  $x$  the \_\_\_\_\_ ( $y$  - values) are approximately normally distributed.
- There is independence in data collection: Data are collected using \_\_\_\_\_ sample of a randomized experiment. When sampling without replacement check that \_\_\_\_\_.

**Checking Conditions** (Be sure to ✓ your conditions!)

To perform a significance test about the slope  $\beta$  of a population regression line, we must check the following:

- The true relationship between  $x$  and  $y$  is linear.

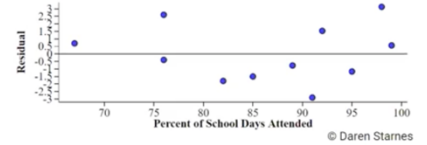
Scatterplot: The scatterplot show a \_\_\_\_\_, \_\_\_\_\_ relationship between the percent of school days attended by a student and the number of correct answers on the state Algebra 1 test.



Residual Plot: The residual plots shows fairly \_\_\_\_\_ scatter about \_\_\_\_\_ line and no obvious curved pattern.

- The standard deviation of  $y$ ,  $\sigma_y$  does not vary with  $x$ .

Scatterplot: The number of correct answers vary by a \_\_\_\_\_ amount for the different percents of school days attended in the data set.



Residual Plot: The residuals are \_\_\_\_\_ in size for the different percents of school days attended in the data set.

- For a particular value of  $x$ , the responses ( $y$  – values) are approximately normally distributed.

Dotplot: A dotplot of residuals shows no obvious \_\_\_\_\_ or \_\_\_\_\_.

- There is independence in data collection.

Data are collected using a random sample or a randomized experiment. The data come from a \_\_\_\_\_ sample of \_\_\_\_\_ Texas Algebra 1 students.

When sampling without replacement, check that  $n \leq 10\%N$ . \_\_\_\_\_ (all Texas Algebra 1 students)

All of the conditions are \_\_\_\_\_.

**What Should We Take Away?**

How do we identify an appropriate significance test procedure for the slope of a population regression line?

When testing a claim about the slope of a population regression line, use a \_\_\_\_\_

How do we verify the conditions for performing a test about the slope of a population regression line?

\* The true relationship between \_\_\_\_\_ is linear. With \_\_\_\_\_ and \_\_\_\_\_

\* The standard deviation of  $y$ ,  $\sigma_y$ , does vary with  $x$ . With \_\_\_\_\_ and \_\_\_\_\_

\* For a particular value of  $x$ , the  $y$ -values are approximately normally distributed. With \_\_\_\_\_ of the \_\_\_\_\_. Need \_\_\_\_\_ if obvious \_\_\_\_\_ or \_\_\_\_\_.

\* There is \_\_\_\_\_ in data collection

Data are collected using a \_\_\_\_\_ sample or a \_\_\_\_\_ experiment.

When sampling \_\_\_\_\_ replacement, check the \_\_\_\_\_.

# AP Statistics CED 9.5 Daily Video 1 (Skill 3.E)

## Carrying Out a Test for the Slope of a Regression Model

### What Will We Learn?

How do we calculate an appropriate test statistic in a test about the slope of a population regression line?

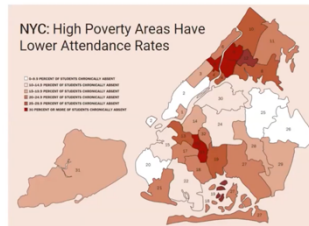
How do we calculate a  $p$ -value in a test about the slope of a population regression line?

### Is Attendance the Solution?

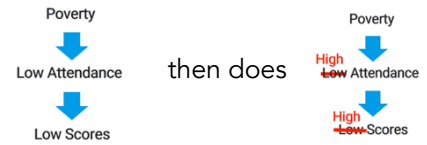
Nationally, higher income areas tend to have fewer chronically absent students.

Possible reasons:

- Transportation access
- Work to support family



Graphic from Neuser et al., "A Better Picture of Poverty," Center for New York City Affairs, Nov. 2014. [https://www.attendanceworks.org/wp-content/uploads/2017/06/BetterPicturePoverty\\_PA\\_FINAL\\_201.pdf](https://www.attendanceworks.org/wp-content/uploads/2017/06/BetterPicturePoverty_PA_FINAL_201.pdf)

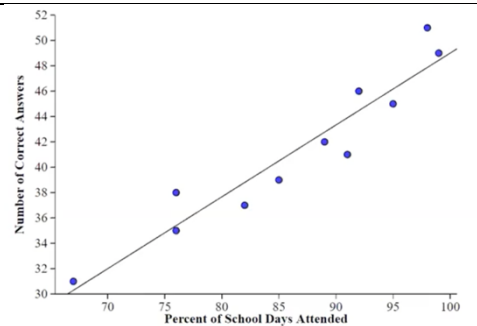


Some school systems have targeted attendance as the key to raising test scores for lower income students.

### Let's Look at the Data

Researchers collected data on the percent of school days attended and the number of questions answered correctly on the state's end-of-course test for a random sample of 11 Texas Algebra 1 students. Here are the data, along with a scatterplot and computer output from a least-squares regression analysis.

Percent attendance	95	89	67	98	99	76	92	91	76	85	82
Questions correct	45	42	31	51	49	38	46	41	35	39	37



Predictor	Coef	SE Coef	T	P
Constant	-7.69	5.37	-1.43	0.186
Attendance	0.57	0.062	9.18	0.000

S = 1.99    R-Sq = 90.3%    R-Sq (adj) = 89.3%

Do the data give convincing evidence at the  $\alpha = 0.01$  significance level of a positive, linear relationship between state test score and percent attendance for Texas Algebra 1 students?

### Test Scores and Attendance – Hypotheses

In a previous video, we stated the hypotheses:  $H_0: \beta = 0$  and  $H_a: \beta > 0$ ; where  $\beta$  = the slope of the population regression line for predicting number of questions answered correctly on the state test from the percent of schools day attended for Texas Algebra 1 students. We will use  $\alpha = 0.01$ . We will use a  $t$ -test for slope and the conditions have all been met.

### Calculating a Test Statistic

In the attendance and test score study,  $b =$  \_\_\_\_\_.

This is evidence for  $H_a: \beta > \underline{\quad}$  because  $b = \underline{\quad} > \underline{\quad}$ .

We want to know how \_\_\_\_\_ it is to get evidence for  $H_a$  this \_\_\_\_\_ by \_\_\_\_\_ alone when  $H_0$  is \_\_\_\_\_.

After verifying the conditions are met, calculate the standardized test statistics:

$$\text{standardized tests statistic} = \frac{\text{statistic} - \text{parameter}}{\text{standard error of the statistic}} \text{ from formula sheet!!}$$

Predictor	Coef	SE Coef	T	P
Constant	-7.69	5.37	-1.43	0.186
Attendance	0.57	0.062	9.18	0.000

S = 1.99    R-Sq = 90.3%    R-Sq (adj) = 89.3%

### Calculating a Test Statistic

For a  $t$ -test for a slope, the standardized test statistic is:  $t = \frac{b - \beta_0}{SE_b}$ ,

where  $\beta_0$  is the value of  $\beta$  specified by the null hypothesis. So, we

would have  $t = \underline{\quad} = \underline{\quad}$

Predictor	Coef	SE Coef	T	P
Constant	-7.69	5.37	-1.43	0.186
Attendance	0.57	0.062	9.18	0.000

S = 1.99    R-Sq = 90.3%    R-Sq (adj) = 89.3%

\* All components can be found on the AP Statistics Formula Sheet! Make sure you can locate them!

**Calculating the p-value**

Once we have calculated the standardized test statistic, use the \_\_\_\_\_ with df = \_\_\_\_\_ to calculate the \_\_\_\_\_. The p-value is the \_\_\_\_\_ of values for the \_\_\_\_\_ distribution that are as \_\_\_\_\_ than the observed value of the test statistic.

Predictor	Coef	SE Coef	T	P
Constant	-7.69	5.37	-1.43	0.186
Attendance	0.57	0.062	9.18	0.000

S = 1.99 R-Sq = 90.3% R-Sq (adj) = 89.3%

Because our alternative hypothesis is  $H_a: \beta > 0$ , we want to find  $P(t \geq \text{_____})$  in a t distribution with df = \_\_\_\_\_.

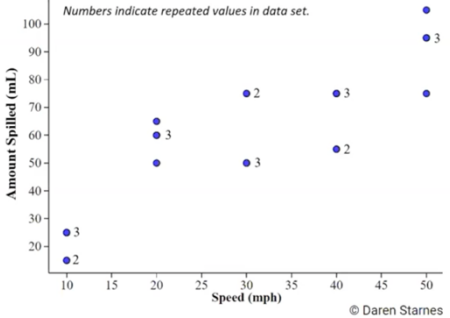
Note: divide by 2 for a one-tailed

You can use Table B or Technology to find the p-value = \_\_\_\_\_ or use the output!

**Don't Spill my Drink!**

Two AP Statistics students wondered if there is a linear relationship between speed and amount of drink spilled when driving on bumpy dirt roads. To find out, they filled a cup with 275 mL of water and placed it in the car's cup holder, then drove down a bumpy dirt road at a specified speed and recorded how much water spilled out. The 25 trials were randomly assigned to be 5 trials each at 10, 20, 30, 40, or 50 miles per hour (mph). Here are the data. (Note that several data points have the same values for speed and amount spilled.) A scatterplot of the data, along with computer output from a least-squares regression analysis are shown. Do the data provide convincing evidence at the 0.05 significance level of a linear relationship between the car's speed on a bumpy dirt road and the amount of drink spilled? Assume the conditions for inference are met. Additionally, we have previously determined the hypotheses.

Speed (mph)	Spilled (mL)	Speed (mph)	Spilled (mL)
10	25	40	75
10	25	40	75
10	25	40	75
10	15	40	55
10	15	40	55
20	60	50	95
20	60	50	95
20	65	50	95
20	60	50	75
20	50	50	105
30	50		
30	75		
30	50		
30	75		
30	50		



Term	Coef	SE Coef
Constant	14.40	5.95
Speed (mph)	1.520	0.179

S=12.6800 R-sq = 75.75%

**Don't Spill My Drink!**

Calculate the standardized test statistic and p-value. From the previous video we know that:  $H_0: \text{_____}$  vs  $H_a: \text{_____}$ ; where  $\beta = \text{_____}$  line for predicting the amount of drink spilled (in mL) from the car's speed (in mph) on a bumpy dirt road. We will use  $\alpha = \text{_____}$  because no significance level was stated. We will conduct a t-test for slope. As a reminder:  $t = \frac{b - \beta_0}{SE_b}$

**Calculating a p-value** (Remember to use Table B or technology!)

Because our alternative hypothesis is  $H_a: \beta \neq 0$ , we want to find  $P(t \leq -8.49) + P(t \geq 8.49)$  in a t distribution with df = \_\_\_\_\_ = \_\_\_\_\_.  
p-value = \_\_\_\_\_

**What Should We Take Away?**

How do we calculate an appropriate test statistic in a test about the slope of a population regression line? **The Formula is  $t = \frac{b - \beta_0}{SE_b}$ . But, often this information can be found on the computer output!**  
How do we calculate a p-value in a test about the slope of a population regression line?  
If \_\_\_\_\_, p-value =  $P(t \geq \text{observed test statistics})$  If \_\_\_\_\_, p-value =  $P(t \geq \text{observed test statistics})$   
If \_\_\_\_\_, p-value =  $2 \times P(t \geq |\text{observed test statistics}|)$ ; in a t distribution with \_\_\_\_\_ df.

# AP Statistics CED 9.5 Daily Video 2 (Skill 4.B)

## Carrying Out a Test for the Slope of a Regression Model

### What Will We Learn?

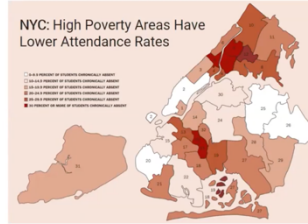
How do we interpret the  $p$ -value in a test about the slope of a population regression line?  
 How do we state a conclusion in a test about the slope of a population regression line?

### Is Attendance the Solution?

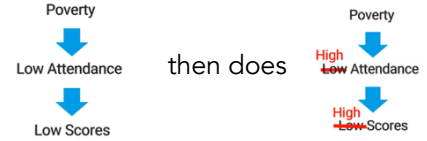
Can education systems equalize opportunities for lower income students?

Nationally, higher income areas tend to have fewer chronically absent students. Possible reasons:

- Transportation access
- Work to support family



Graphic from Nasar et al. "A Better Picture of Poverty," Center for New York City Affairs, Nov. 2014. [http://www.attendancerevival.org/wp-content/uploads/2017/08/BetterPictureofPoverty\\_Full\\_PDF\\_001.pdf](http://www.attendancerevival.org/wp-content/uploads/2017/08/BetterPictureofPoverty_Full_PDF_001.pdf)

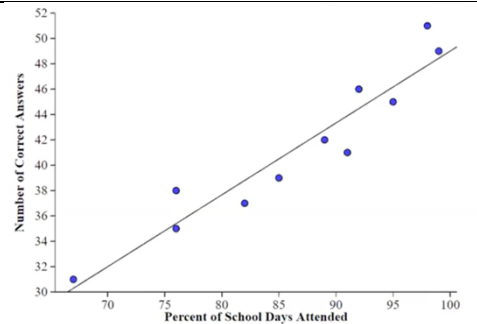


Some school systems have targeted attendance as the key to raising test scores for lower income students.

### Let's Look at the Data

Researchers collected data on the percent of school days attended and the number of questions answered correctly on the state's end-of-course test for a random sample of 11 Texas Algebra 1 students. Here are the data, along with a scatterplot and computer output from a least-squares regression analysis.

Percent attendance	95	89	67	98	99	76	92	91	76	85	82
Questions correct	45	42	31	51	49	38	46	41	35	39	37



© Daren Starnes

Predictor	Coef	SE Coef	T	P
Constant	-7.69	5.37	-1.43	0.186
Attendance	0.57	0.062	9.18	0.000

S = 1.99    R-Sq = 90.3%    R-Sq (adj) = 89.3%

Do the data give convincing evidence at the  $\alpha = 0.01$  significance level of a positive, linear relationship between state test score and percent attendance for Texas Algebra 1 students?

**Test Scores and Attendance - Hypotheses** In a previous video, we stated the hypotheses:  $H_0: \beta = 0$  and  $H_a: \beta > 0$ ; where  $\beta$  = the slope of the population regression line for predicting number of questions answered correctly on the state test from the percent of schools day attended for Texas Algebra 1 students. Use  $\alpha = 0.01$ . We will use a t-test for slope and the conditions have all been met.

**Interpreting p-value** - From previous videos:

$b =$  \_\_\_\_\_,  $t =$  \_\_\_\_\_ and  $p$ -value  $\approx$  \_\_\_\_\_  
 $b =$  \_\_\_\_\_ is evidence \_\_\_\_\_  $H_0: \beta =$  \_\_\_\_\_ and for  $H_a: \beta >$  \_\_\_\_\_, because \_\_\_\_\_.  
 The  $p$ -value measures how \_\_\_\_\_ it is to get evidence for  $H_a$  as \_\_\_\_\_ than the observed evidence by \_\_\_\_\_ alone when  $H_0$  is \_\_\_\_\_. Assuming there is \_\_\_\_\_ relationship between state test score and percent attendance for Texas Algebra 1 students, there is an \_\_\_\_\_ probability of getting a sample regression line with slope \_\_\_\_\_ or \_\_\_\_\_ by \_\_\_\_\_ alone in a random sample of 11 Texas Algebra 1 students.

### Making a Conclusion

Small  $p$ -values  $\longrightarrow$  test statistic \_\_\_\_\_ to occur by \_\_\_\_\_ chance alone.  
 Large  $p$ -values  $\longrightarrow$  test statistic \_\_\_\_\_ to occur by \_\_\_\_\_ chance alone.

- Because the  $p$ -value of \_\_\_\_\_  $\leq \alpha =$  \_\_\_\_\_, we reject  $H_0$ .  
 There is convincing \_\_\_\_\_ evidence that [\_\_\_\_\_].
- Because the  $p$ -value of \_\_\_\_\_  $\leq \alpha =$  \_\_\_\_\_, we fail to reject  $H_0$ .  
 There is not convincing \_\_\_\_\_ evidence that [\_\_\_\_\_].



**Making a Conclusion**

Do the data give convincing evidence at the  $\alpha = 0.01$  significance level of a positive, linear relationship between state test score and percent attendance for Texas Algebra 1 students?

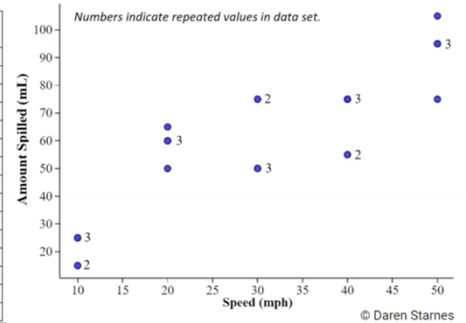
- $H_0: \beta$  \_\_\_\_\_ vs.  $H_a: \beta$  \_\_\_\_\_, where  $\beta$  = the slope of the population regression line for predicting number of questions answered correctly on the state test from the percent of school days attended for Texas Algebra 1 students.
- Conditions are \_\_\_\_\_
- $b =$  \_\_\_\_\_,  $t =$  \_\_\_\_\_ and  $p$ -value  $\approx$  \_\_\_\_\_

Because the  $p$ -value of \_\_\_\_\_  $<$  \_\_\_\_\_, we reject  $H_0$ . There \_\_\_\_\_ convincing statistical evidence of a \_\_\_\_\_ between state test score and percent attendance for Texas Algebra 1 students.

**Don't Spill my Drink!**

Two AP Statistics students wondered if there is a linear relationship between speed and amount of drink spilled when driving on bumpy dirt roads. To find out, they filled a cup with 275 mL of water and placed it in the car's cup holder, then drove down a bumpy road at a specified speed and recorded how much water spilled out. The 25 trials were randomly assigned to be 5 trials each at 10, 20, 30, 40, or 50 miles per hour (mph).

Speed (mph)	Spilled (mL)	Speed (mph)	Spilled (mL)
10	25	40	75
10	25	40	75
10	25	40	75
10	15	40	55
10	15	40	55
20	60	50	95
20	60	50	95
20	65	50	95
20	60	50	75
20	50	50	105
30	50		
30	75		
30	50		
30	75		
30	50		



Term	Coef	SE Coef
Constant	14.40	5.95
Speed (mph)	1.520	0.179
S=12.6800		R-sq = 75.75%

Here are the data. (Note that several data points have the same values for speed and amount spilled.) A scatterplot of the data, along with computer output from a least-squares regression analysis are shown. Do the data provide convincing evidence at the 0.05 significance level of a linear relationship between the car's speed on a bumpy dirt road and the amount of drink spilled? Assume the conditions for inference are met.

**Interpreting the p-value**

From previous videos:

- $H_0: \beta$  \_\_\_\_\_ vs.  $H_a: \beta$  \_\_\_\_\_, where  $\beta$  = the slope of the \_\_\_\_\_ regression line for predicting amount of \_\_\_\_\_ spilled (in mL) from the car's speed (in mph) on a bumpy dirt road.
- $t$ -test for slope: Conditions are met.
- $b =$  \_\_\_\_\_,  $t =$  \_\_\_\_\_ and  $p$ -value  $\approx$  \_\_\_\_\_

Interpret the  $p$ -value.

Assuming that there is \_\_\_\_\_ relationship between the car's speed on a bumpy dirt road and the amount of drink spilled, there is a \_\_\_\_\_ probability of getting a sample regression line with a slope as \_\_\_\_\_ than \_\_\_\_\_ in either direction by \_\_\_\_\_ alone.

**Making a Conclusion**

Do the data provide convincing evidence at the 0.05 significance level of a linear relationship between the car's speed on a bumpy dirt road and the amount of drink spilled? From previous video:

- $H_0: \beta$  \_\_\_\_\_ vs.  $H_a: \beta$  \_\_\_\_\_, where  $\beta$  = the slope of the \_\_\_\_\_ regression line for predicting amount of \_\_\_\_\_ spilled (in mL) from the car's speed (in mph) on a bumpy dirt road.
- Conditions are met. And  $b =$  \_\_\_\_\_,  $t =$  \_\_\_\_\_ and  $p$ -value  $\approx$  \_\_\_\_\_

Because the  $p$ -value of \_\_\_\_\_  $<$  \_\_\_\_\_, we \_\_\_\_\_. There is convincing statistical evidence of a \_\_\_\_\_ relationship between the car's speed on a bumpy dirt road and the amount of drink spilled.

**What Should We Take Away?**

How do we interpret the  $p$ -value in a test about the slope of a population regression line?

The  $p$ -value measures how \_\_\_\_\_ it is to get evidence for  $H_a$  as \_\_\_\_\_ or \_\_\_\_\_ than the observed evidence by chance along with  $H_0$  is \_\_\_\_\_.

How do we state a conclusion in a test about the slope of a population regression line?

- Because the  $p$ -value of \_\_\_\_\_  $\leq \alpha =$  \_\_\_\_\_, we reject  $H_0$ .  
There is convincing \_\_\_\_\_ evidence that [\_\_\_\_\_].
- Because the  $p$ -value of \_\_\_\_\_  $\leq \alpha =$  \_\_\_\_\_, we fail to reject  $H_0$ .  
There is not convincing \_\_\_\_\_ evidence that [\_\_\_\_\_].

## AP Statistics CED 9.5 Daily Video 3 (Skill 4.B)

### Carrying Out a Test for the Slope of a Regression Model

#### What Will We Learn?

How do we perform a complete significance test about the slope of a population regression line?

#### 2001 Exam #6 (Modified)

The Statistics Department at a large university is trying to determine if it is possible to predict whether an applicant will successfully complete the Ph.D. program or will leave before completing the program. The department is considering whether GPA (grade point average) in undergraduate statistics and mathematics courses (a measure of performance) and mean number of credit hours per semester (a measure of workload) would be helpful measures. To gather data, a random sample of 20 entering students from the past 5 years is taken. The data are given below:

Successfully Completed Ph.D. Program

Student	A	B	C	D	E	F	G	H	I	J	K	L	M
GPA	3.8	3.5	4.0	3.9	2.9	3.5	3.5	4.0	3.9	3.0	3.4	3.7	3.6
Credit hours	12.7	13.1	12.5	13.0	15.0	14.7	14.5	12.0	13.1	15.3	14.6	12.5	14.0

Did Not Complete Ph.D. Program

Student	N	O	P	Q	R	S	T
GPA	3.6	2.9	3.1	3.5	3.9	3.6	3.3
Credit hours	11.1	14.5	14.0	10.9	11.5	12.1	12.0

The regression output below resulted from fitting a line to the data in the group of 7 students that did not complete the Ph.D. program. The residual plot (not shown) indicated no unusual patterns, and the assumptions necessary for inference were judged to be reasonable.

Did Not Complete Ph.D. Program

Predictor	Coef	StDev	T	P
Constant	24.200	3.474	6.97	0.001
GPA	-3.485	1.013	-3.44	0.018
S = 0.8408		R-Sq = 70.3%		

(a) For the students who did not complete the Ph.D. program, is there a significant relationship between GPA and mean number of credit hours per semester at the  $\alpha = 0.01$  level?

$H_0: \beta = 0$  and  $H_a: \beta \neq 0$ ; where  $\beta =$  the slope of the \_\_\_\_\_ regression line for predicting \_\_\_\_\_ number of credit hours per semester from \_\_\_\_\_ for students who did \_\_\_\_\_ complete the statistics Ph.D. program at this large university. Use  $\alpha = \underline{\quad}$ ;  $t$ -test for slope; Conditions are met!

#### 2001 Exam #6 (Modified)

(a) For the students who did not complete the Ph.D. program, is there a significant relationship between GPA and mean number of credit hours per semester at the  $\alpha = 0.01$  level?

Did Not Complete Ph.D. Program

Predictor	Coef	StDev	T	P
Constant	24.200	3.474	6.97	0.001
GPA	-3.485	1.013	-3.44	0.018
S = 0.8408		R-Sq = 70.3%		

$t = \underline{\quad}$ ;  $df = \underline{\quad}$ ;  $p$ -value = \_\_\_\_\_

Because the  $p$ -value of \_\_\_\_\_, we \_\_\_\_\_. There is \_\_\_\_\_ convincing evidence of a \_\_\_\_\_ relationship between GPA and mean number of credit hours in the population of students at this university who did not complete the statistics Ph.D. program.

#### 2001 Exam #6 (Modified)

The regression output below resulted from fitting a line to the data in the group of 13 students that successfully completed the Ph.D. program.

The residual plot (not shown) indicated no unusual patterns, and the assumptions necessary for inference were judged to be reasonable.

Successfully Completed Ph.D. Program

Predictor	Coef	StDev	T	P
Constant	23.514	1.684	13.95	0.000
GPA	-2.7555	0.4668	-5.90	0.000
S = 0.5658		R-Sq = 76.0%		

(b) For the students who successfully completed the Ph.D. program, is the evidence for a significant (linear) relationship between GPA and mean number of credit hours per semester stronger or weaker than for the students who did not complete the Ph.D. program? Justify your answer.

Name \_\_\_\_\_

### 2001 Exam #6 (Modified)

Hint: You do not need to run a full test to answer this question. You only need to use the computer output to the right.

Successfully Completed Ph.D. Program

Predictor	Coef	StDev	T	P
Constant	23.514	1.684	13.95	0.000
GPA	-2.7555	0.4668	-5.90	0.000
S = 0.5658		R-Sq = 76.0%		

\_\_\_\_\_. The  $p$ -value for a test of  $H_0: \beta = 0$  and  $H_a: \beta \neq 0$ ; where  $\beta$  = the \_\_\_\_\_ of the population regression line for \_\_\_\_\_ number of credit hours per semester \_\_\_\_\_ GPA for students who successfully completed the statistics Ph.D. program at this large university, is \_\_\_\_\_. This  $p$ -value is \_\_\_\_\_ than the \_\_\_\_\_  $p$ -value for the test in part (a), and gives \_\_\_\_\_ evidence of a \_\_\_\_\_ relationship between GPA and the mean number of credit hours per semester. You could also compare the \_\_\_\_\_ from the tests!

### What Should We Take Away?

How do we perform a complete significance test about the slope of a population regression line?

Make sure to:

- State the \_\_\_\_\_ and \_\_\_\_\_ hypotheses, and define \_\_\_\_\_.
- Give the \_\_\_\_\_.
- Identify the \_\_\_\_\_ you are using.
- Verify that the \_\_\_\_\_ for the procedure are \_\_\_\_\_.
- Calculate the \_\_\_\_\_ and the \_\_\_\_\_.
- Make a \_\_\_\_\_ based on the \_\_\_\_\_. (You do \_\_\_\_\_ need to interpret the  $p$ -value unless specifically asked.)

**AP Statistics CED 9.6 Daily Video 1****Carrying Out a Test for the Slope of a Regression Model****What Will We Learn?**

What inference methods did we learn about in Units 6 – 9?

How can we identify the appropriate inference procedure to use in a given setting?

**Inference Recap**

Two main goals of inference:

Estimating a parameter: \_\_\_\_\_

Testing a claim: \_\_\_\_\_

Appropriate inference method depends on type of data:

Unit 6 Inference for Categorical Data: \_\_\_\_\_

Unit 7 Inference for Quantitative Data: \_\_\_\_\_

Unit 8 Inference for Categorical Data: \_\_\_\_\_

Unit 9 Inference for Categorical Data: \_\_\_\_\_

**Confidence Intervals**

The formula on the formula sheet is the same for \_\_\_\_\_ procedures.

$$CI = \text{_____} \pm (\text{_____})(\text{_____})$$

You should always do these steps:

- Define the \_\_\_\_\_ you are trying to estimate.
- Identify the inference \_\_\_\_\_.
- Verify the \_\_\_\_\_ for the procedure are \_\_\_\_\_.
- \_\_\_\_\_ the confidence interval.
- \_\_\_\_\_ the interval \_\_\_\_\_.

You do \_\_\_\_\_ need to interpret the confidence \_\_\_\_\_ unless specifically asked.

**Significance Tests**

$$\textit{standardized test statistic} = \frac{\textit{statistic} - \textit{parameter}}{\textit{standard error of the statistic}}$$

You should always do these steps:

- State the \_\_\_\_\_ and \_\_\_\_\_ hypotheses. Be sure to \_\_\_\_\_ parameter.
- Give the \_\_\_\_\_.
- \_\_\_\_\_ the inference procedure.
- Verify that the \_\_\_\_\_ for the procedure are \_\_\_\_\_.
- Calculate the \_\_\_\_\_ and the \_\_\_\_\_.
- Make a conclusion based on the \_\_\_\_\_.

You do \_\_\_\_\_ need to interpret the \_\_\_\_\_ unless specifically asked.

**Unit 6**

Inference for Categorical Data: Proportions			
Goal	Method	Formula	Conditions
Estimate $p$	One-sample z interval for a proportion	$\hat{p} \pm z^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$	<ul style="list-style-type: none"> <li>Random sample/random assignment</li> <li><math>n \leq 10\%N</math> (samp w/o replacement)</li> <li><math>n\hat{p}</math> and <math>n(1-\hat{p}) \geq 10</math></li> </ul>
Test $H_0: p = p_0$	One-sample z test for a proportion	$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$	<ul style="list-style-type: none"> <li>Random sample/random assignment</li> <li><math>n \leq 10\%N</math> (samp w/o replacement)</li> <li><math>np_0 \geq 10</math> and <math>n(1-p_0) \geq 10</math></li> </ul>
Estimate $p_1 - p_2$	Two-sample z interval for difference in proportions	$(\hat{p}_1 - \hat{p}_2) \pm z^* \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$	<ul style="list-style-type: none"> <li>Random samples/random assignment</li> <li><math>n_1 \leq 10\%N_1</math> and <math>n_2 \leq 10\%N_2</math> (samp w/o replacement)</li> <li><math>n_1\hat{p}_1, n_1(1-\hat{p}_1), n_2\hat{p}_2, n_2(1-\hat{p}_2) \geq 10</math></li> </ul>
Test $H_0: p_1 - p_2 = 0$	Two-sample z test for a difference in proportions	$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\frac{\hat{p}_c(1-\hat{p}_c)}{n_1} + \frac{\hat{p}_c(1-\hat{p}_c)}{n_2}}}$	<ul style="list-style-type: none"> <li>Random samples/random assignment</li> <li><math>n_1 \leq 10\%N_1</math> and <math>n_2 \leq 10\%N_2</math> (samp w/o replacement)</li> <li><math>n_1\hat{p}_c, n_1(1-\hat{p}_c), n_2\hat{p}_c, n_2(1-\hat{p}_c) \geq 10</math></li> </ul>

**Unit 7**

Inference for Quantitative Data: Means			
Goal	Method	Formula	Conditions
Estimate $\mu$	One-sample t interval for a mean	$\bar{x} \pm t^* \frac{s}{\sqrt{n}}$	<ul style="list-style-type: none"> <li>Random sample/random assignment</li> <li><math>n \leq 10\%N</math> (samp w/o replacement)</li> <li><math>n \geq 30</math> or no strong skew/outliers</li> </ul>
Test $H_0: \mu = \mu_0$	One-sample t test for a mean	$t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$	<ul style="list-style-type: none"> <li>Random sample/random assignment</li> <li><math>n \leq 10\%N</math> (samp w/o replacement)</li> <li><math>n \geq 30</math> or no strong skew/outliers</li> </ul>
Estimate $\mu_1 - \mu_2$	Two-sample t interval for difference in means	$(\bar{x}_1 - \bar{x}_2) \pm t^* \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$	<ul style="list-style-type: none"> <li>Random samples/random assignment</li> <li><math>n_1 \leq 10\%N_1</math> and <math>n_2 \leq 10\%N_2</math> (samp w/o replacement)</li> <li><math>n_1, n_2 \geq 30</math> or no strong skew/outliers</li> </ul>
Test $H_0: \mu_1 - \mu_2 = 0$	Two-sample t test for difference in means	$t = \frac{(\bar{x}_1 - \bar{x}_2) - 0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$	<ul style="list-style-type: none"> <li>Random samples/random assignment</li> <li><math>n_1 \leq 10\%N_1</math> and <math>n_2 \leq 10\%N_2</math> (samp w/o replacement)</li> <li><math>n_1, n_2 \geq 30</math> or no strong skew/outliers</li> </ul>

**Unit 8**

Inference for Categorical Data: Chi-Square			
Goal	Method	Formula	Conditions
Test $H_0$ : Categorical variable has specified distribution	Chi-square test for goodness of fit	$\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$	<ul style="list-style-type: none"> <li>Random sample/random assignment</li> <li><math>n \leq 10\%N</math> (samp w/o replacement)</li> <li>All expected counts &gt; 5.</li> </ul>
Test $H_0$ : Categorical variable has same distribution for each population or treatment	Chi-square test for homogeneity	$\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$	<ul style="list-style-type: none"> <li>Random samples/random assignment</li> <li><math>n_1 \leq 10\%N_1, n_2 \leq 10\%N_2, \dots</math> (samp w/o replacement)</li> <li>All expected counts &gt; 5.</li> </ul>
Test $H_0$ : There is no association between two categorical variables in a population	Chi-square test for independence	$\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$	<ul style="list-style-type: none"> <li>Random sample/random assignment</li> <li><math>n \leq 10\%N</math> (samp w/o replacement)</li> <li>All expected counts &gt; 5.</li> </ul>



Inference for Quantitative Data: Slopes			
Goal	Method	Formula	Conditions
Estimate $\beta$	$t$ interval for a slope	$b \pm t^*SE_b = b \pm t^* \frac{s}{s_x \sqrt{n-1}}$	<ul style="list-style-type: none"> <li>The true relationship between <math>x</math> and <math>y</math> is linear.</li> <li>The standard deviation of <math>y</math> doesn't vary with <math>x</math>.</li> <li>For a particular value of <math>x</math>, the <math>y</math> values are approximately normally distributed.</li> <li>There is independence in data collection.                             <ul style="list-style-type: none"> <li>Random sample or randomized experiment.</li> <li><math>n \leq 10\%N</math> (samp w/o replacement)</li> </ul> </li> </ul>
Test $H_0: \beta = \beta_0$	$t$ test for a slope	$t = \frac{b - \beta_0}{SE_b}$	<ul style="list-style-type: none"> <li>The true relationship between <math>x</math> and <math>y</math> is linear.</li> <li>The standard deviation of <math>y</math> doesn't vary with <math>x</math>.</li> <li>For a particular value of <math>x</math>, the <math>y</math> values are approximately normally distributed.</li> <li>There is independence in data collection.                             <ul style="list-style-type: none"> <li>Random sample or randomized experiment.</li> <li><math>n \leq 10\%N</math> (samp w/o replacement)</li> </ul> </li> </ul>

**Selecting an Inference Procedure**

Inference for ...			
Categorical Data: Proportions	Quantitative Data: Means	Categorical Data: Chi-Square	Quantitative Data: Slopes
One-sample $z$ interval for a proportion	One-sample $t$ interval for a mean (Paired data)	Chi-square test for goodness of fit (Distribution of proportions for one categorical variable)	$t$ interval for a slope
One-sample $z$ test for a proportion	One-sample $t$ test for a mean (Paired data)	Chi-square test for homogeneity (Distribution of a categorical variable for multiple populations or treatments)	$t$ test for a slope
Two-sample $z$ interval for a difference in proportions	Two-sample $t$ interval for a difference in means	Chi-square test for independence (Relationship between two categorical variables)	
Two-sample $z$ test for a difference in proportions	Two-sample $t$ test for a difference in means		

**1997 Exam #5**

A company bakes computer chips in two ovens, oven A and oven B. The chips are randomly assigned to an oven and hundreds of chips are baked each hour. The percentage of defective chips coming from these ovens for each hour of production throughout a day is shown below.

The percentage of defective chips produced each hour by oven A has a mean of 33.56 and a standard deviation of 5.20. The percentage of defective chips produced each hour by oven B has a mean of 32.44 and a standard deviation of 3.78. The hourly differences in percentages for oven A minus oven B have a mean of 1.11 and a standard deviation of 4.28.

Hour	Oven A	Oven B
1	45	36
2	32	37
3	34	33
4	31	34
5	35	33
6	37	32
7	31	33
8	30	30
9	27	24

Does there appear to be a difference between oven A and oven B with respect to the mean percentages of defective chips produced? Give appropriate statistical evidence to support your answer.

Individual/case: \_\_\_\_\_ Inference for: \_\_\_\_\_  
 Variable(s) of interest: \_\_\_\_\_ One sample, Paired data, or Two Samples?  
 Type of data: \_\_\_\_\_ Estimate parameter or Test a Claim? Perform: \_\_\_\_\_

**Larry Green's Applet**

Source: Larry Green's Comprehensive Review applet,  
<http://www.Itconline.net/green/java/Statistics/catStatProb/categorizingStatProblems.JavaScript.html>

Is political affiliation related to the month that the person was born in? 3000 voters were studied.

Individual/case: \_\_\_\_\_ Variable(s) of interest: \_\_\_\_\_

Type of Data: \_\_\_\_\_ Inference for: \_\_\_\_\_

1 Variable, 2 Variables (one sample) or 2 Variables (multiple samples or treatments)?

Inference Procedure: \_\_\_\_\_

- Confidence Interval for a Population Mean
- Confidence Interval for a Proportion
- Confidence Interval for the Diff. Between 2 Means (Independent Samples)
- Confidence Interval for Paired Data (Dependent Samples)
- Confidence Interval for the Difference Between 2 Proportions
- Prediction for a Single Value of  $y$  for a Fixed  $x$
- Hypothesis Test for a Population Mean
- Hypothesis Test for a Population Proportion
- Hyp. Test for the Difference Between 2 Means (Independent Samples)
- Hyp. Test for Paired Data (Dependent Samples)
- Hyp. Test for the Difference Between 2 Proportions
- Chi-Square Goodness of Fit Test
- Chi-Square Test for Independence
- Chi-Square Test for Homogeneity

Is political affiliation related to the month that the person was born in? 3000 voters were studied.

Individual/case: \_\_\_\_\_ Variable(s) of interest: \_\_\_\_\_

Type of Data: \_\_\_\_\_ Inference for: \_\_\_\_\_

1 Variable, 2 Variables (one sample) or 2 Variables (multiple samples or treatments)?

Inference Procedure: \_\_\_\_\_

**What Should We Take Away?**

What inference methods did we learn about in Units 6 – 9?

**Unit 6 Inference for Categorical Data:** \_\_\_\_\_

**Unit 7 Inference for Quantitative Data:** \_\_\_\_\_

**Unit 8 Inference for Categorical Data:** \_\_\_\_\_

**Unit 9 Inference for Categorical Data:** \_\_\_\_\_

How can we identify the appropriate inference procedure to use in a given setting?

**Individual/case:** \_\_\_\_\_

**Variables of interest:** \_\_\_\_\_

**Type of Data:** \_\_\_\_\_

**Inference for:** \_\_\_\_\_