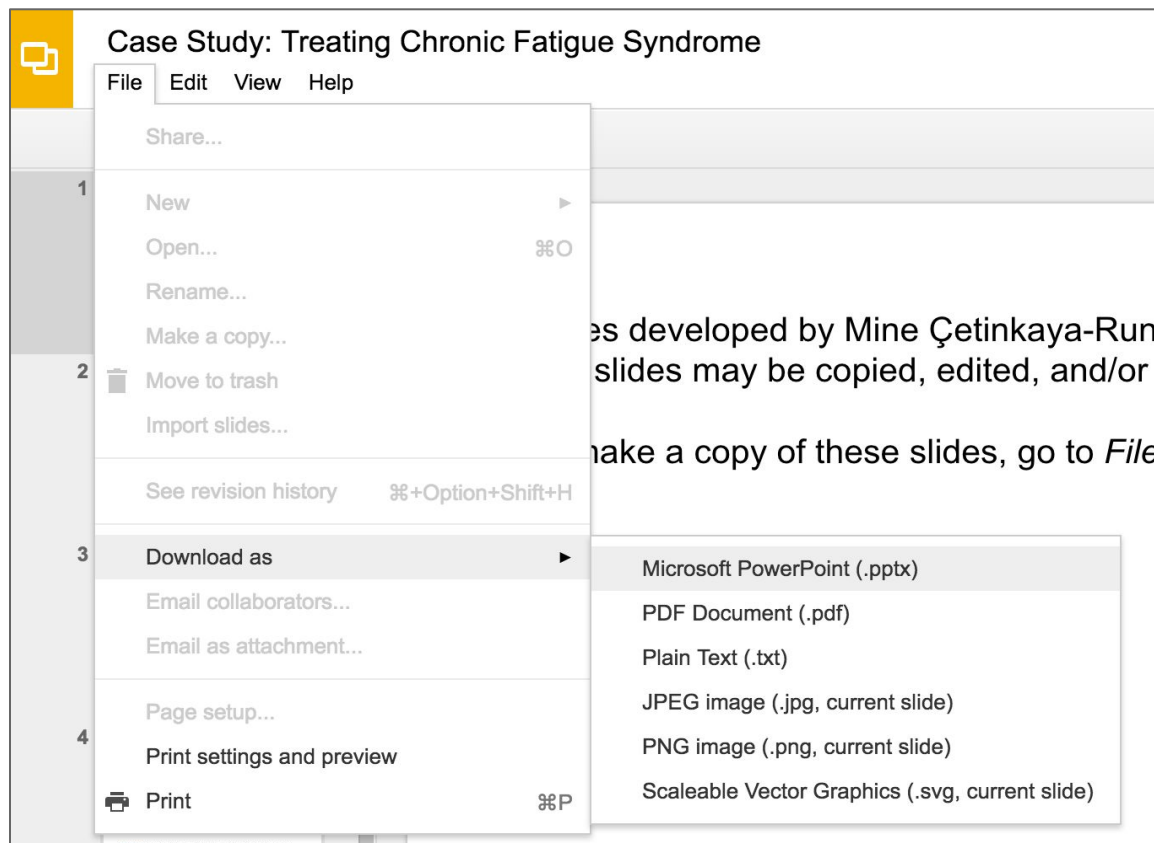


Slides developed by Mine Çetinkaya-Rundel of OpenIntro  
Translated from LaTeX to Google Slides by Curry W. Hilton of OpenIntro.  
The slides may be copied, edited, and/or shared via the [CC BY-SA license](https://creativecommons.org/licenses/by-sa/4.0/)

To make a copy of these slides, go to *File > Download as > [option]*, as shown below. Or if you are logged into a Google account, you can choose *Make a copy...* to create your own version in Google Drive.



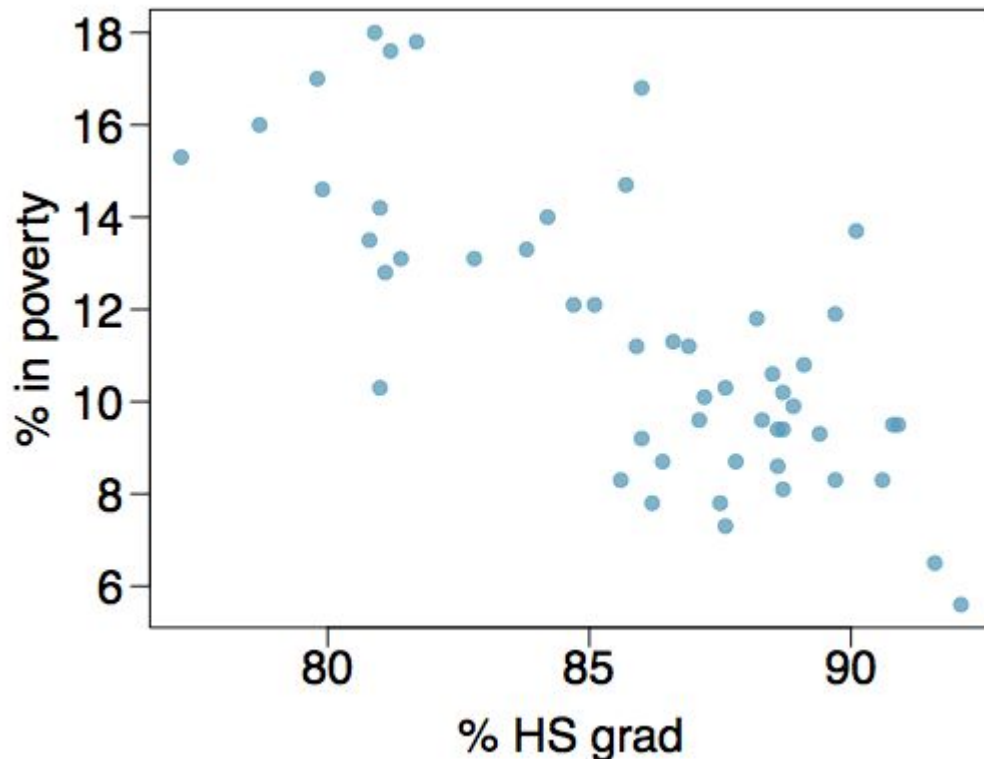
# **Line Fitting, Residuals, and Correlation**

# Modeling numerical variables

In this unit we will learn to quantify the relationship between two numerical variables, as well as modeling numerical response variables using a numerical or categorical explanatory variable.

# Poverty vs. HS graduate rate

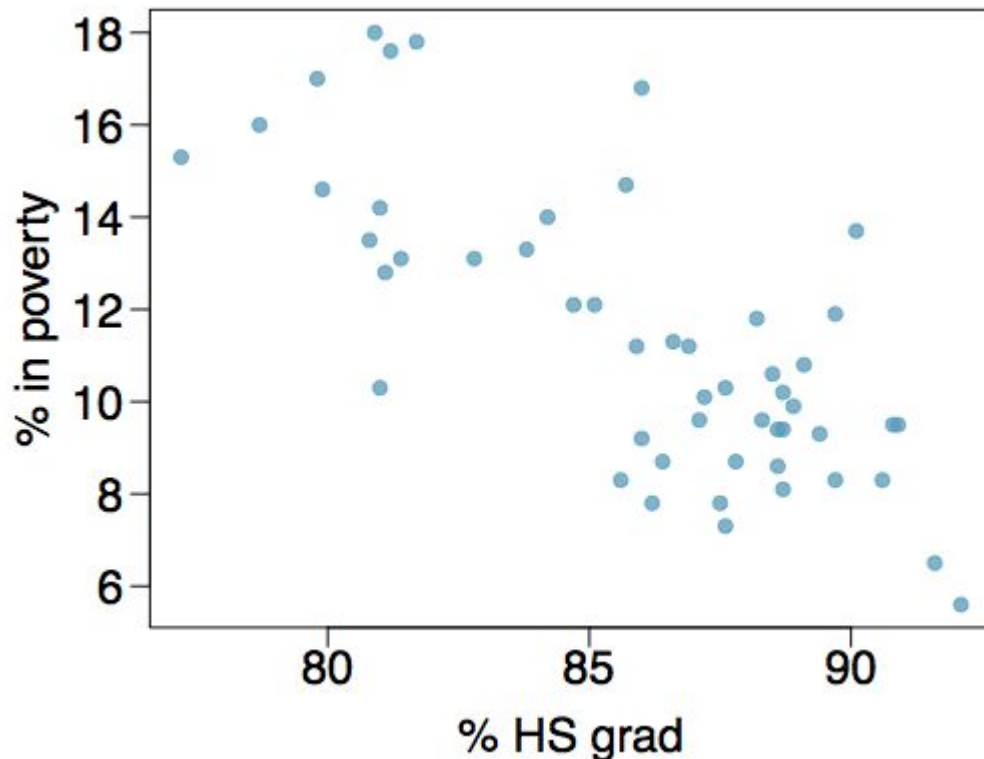
The *scatterplot* below shows the relationship between HS graduate rate in all 50 US states and DC and the percent of residents who live below the poverty line (income below \$23,050 for a family of 4 in 2012).



Response variable?

# Poverty vs. HS graduate rate

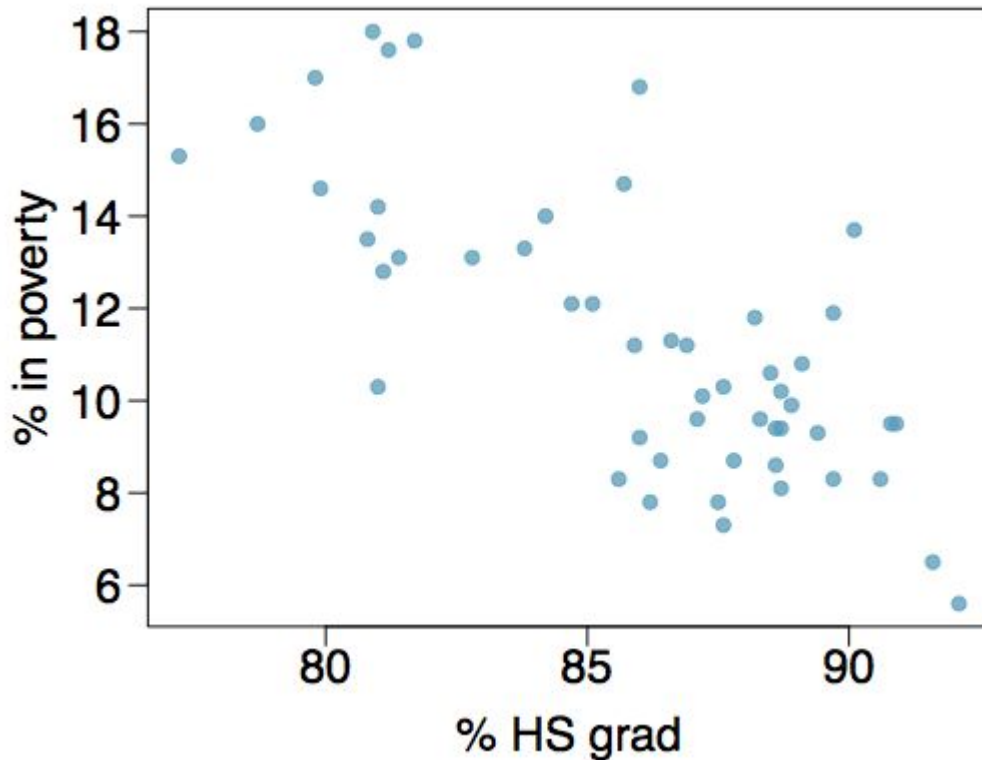
The *scatterplot* below shows the relationship between HS graduate rate in all 50 US states and DC and the percent of residents who live below the poverty line (income below \$23,050 for a family of 4 in 2012).



Response variable?  
*% in poverty*

# Poverty vs. HS graduate rate

The *scatterplot* below shows the relationship between HS graduate rate in all 50 US states and DC and the percent of residents who live below the poverty line (income below \$23,050 for a family of 4 in 2012).

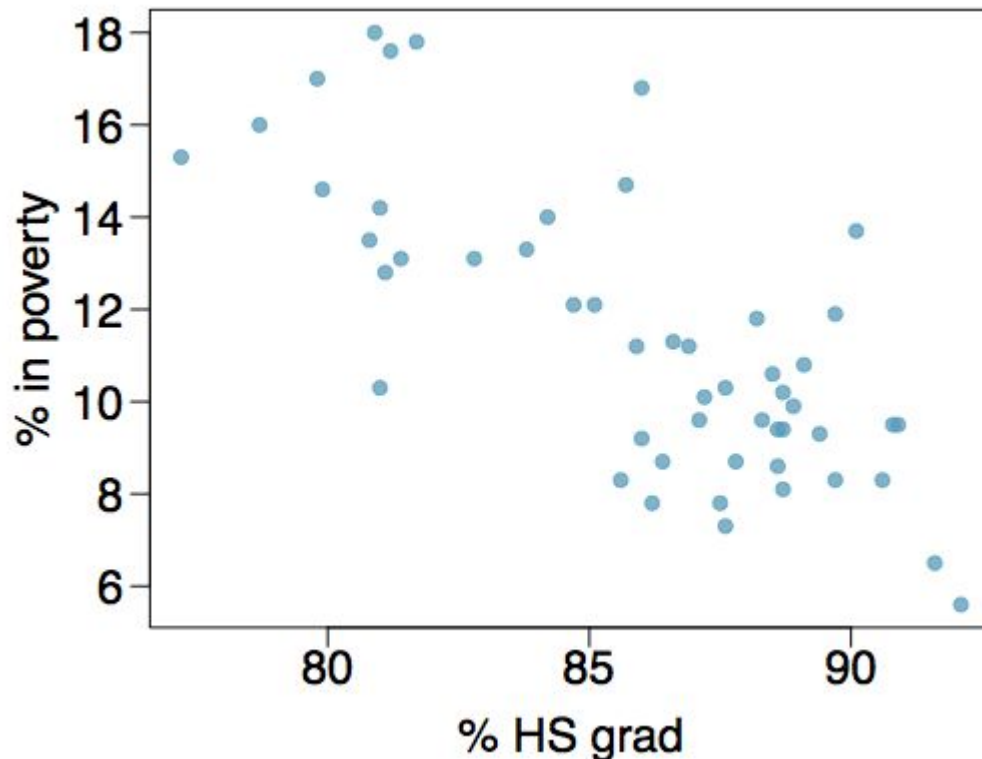


Response variable?  
*% in poverty*

## Explanatory variable?

# Poverty vs. HS graduate rate

The *scatterplot* below shows the relationship between HS graduate rate in all 50 US states and DC and the percent of residents who live below the poverty line (income below \$23,050 for a family of 4 in 2012).



Response variable?

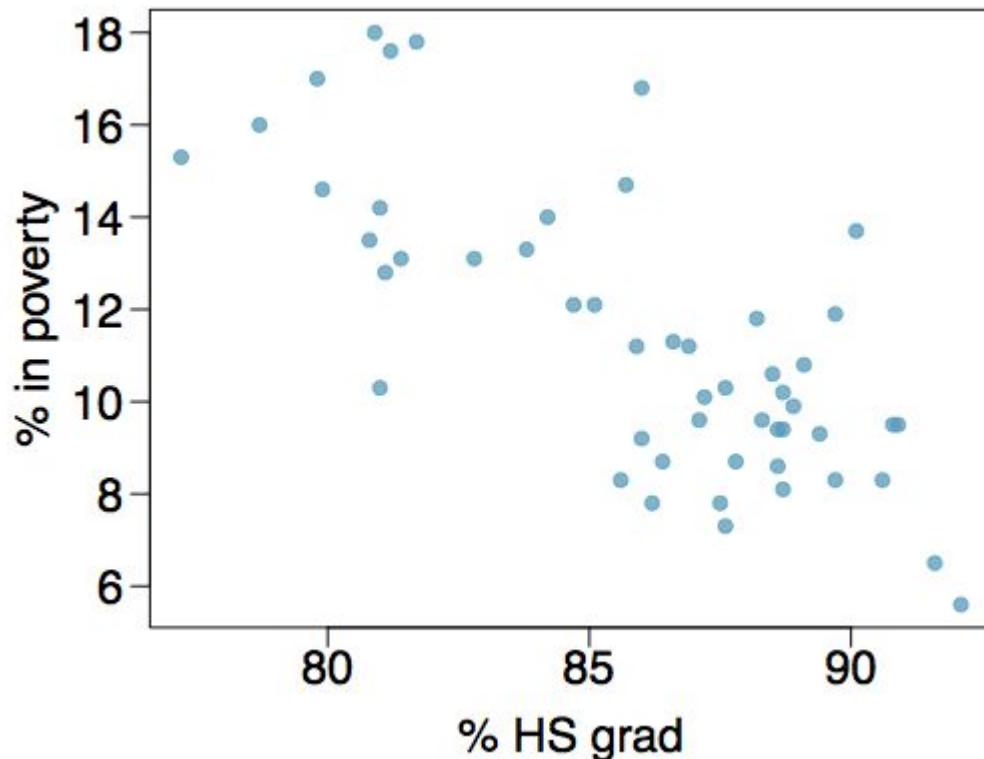
*% in poverty*

Explanatory variable?

*% HS grad*

# Poverty vs. HS graduate rate

The *scatterplot* below shows the relationship between HS graduate rate in all 50 US states and DC and the percent of residents who live below the poverty line (income below \$23,050 for a family of 4 in 2012).



Response variable?

*% in poverty*

Explanatory variable?

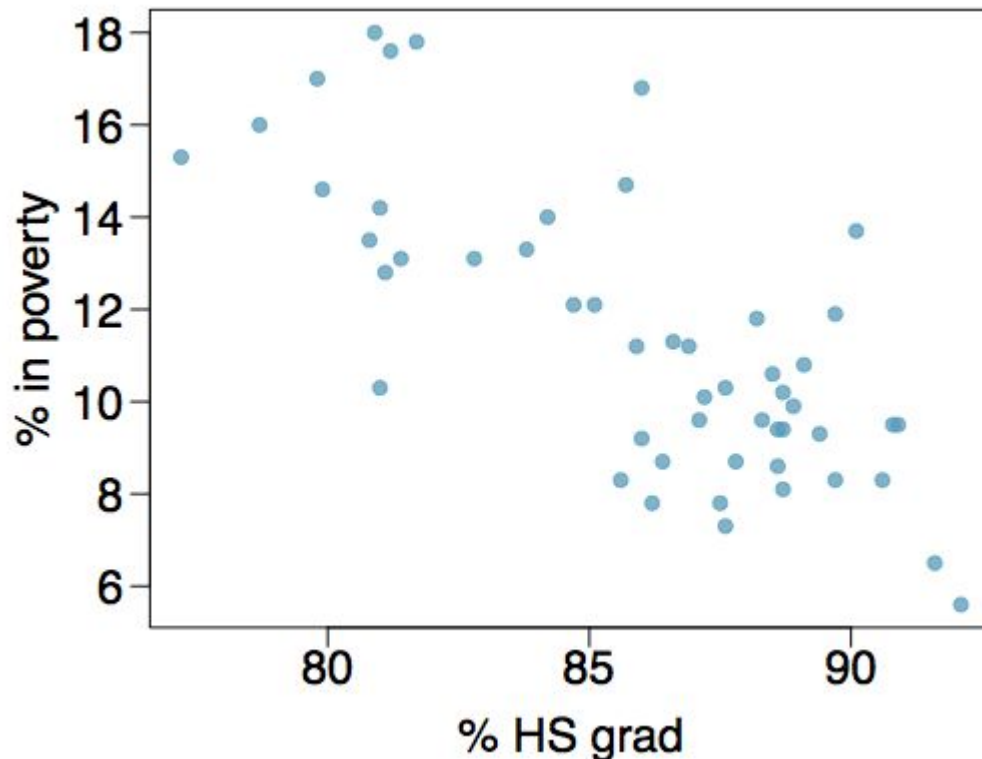
*% HS grad*

Relationship?



# Poverty vs. HS graduate rate

The *scatterplot* below shows the relationship between HS graduate rate in all 50 US states and DC and the percent of residents who live below the poverty line (income below \$23,050 for a family of 4 in 2012).



Response variable?

*% in poverty*

Explanatory variable?

*% HS grad*

Relationship?

*linear, negative,  
moderately strong*

# Poverty vs. HS graduate rate

The linear model for predicting poverty from high school graduation rate in the US is

$$\hat{poverty} = 64.78 - 0.62 * HS_{grad}$$

The "hat" is used to signify that this is an estimate.

# Poverty vs. HS graduate rate

The linear model for predicting poverty from high school graduation rate in the US is

$$\hat{poverty} = 64.78 - 0.62 * HS_{grad}$$

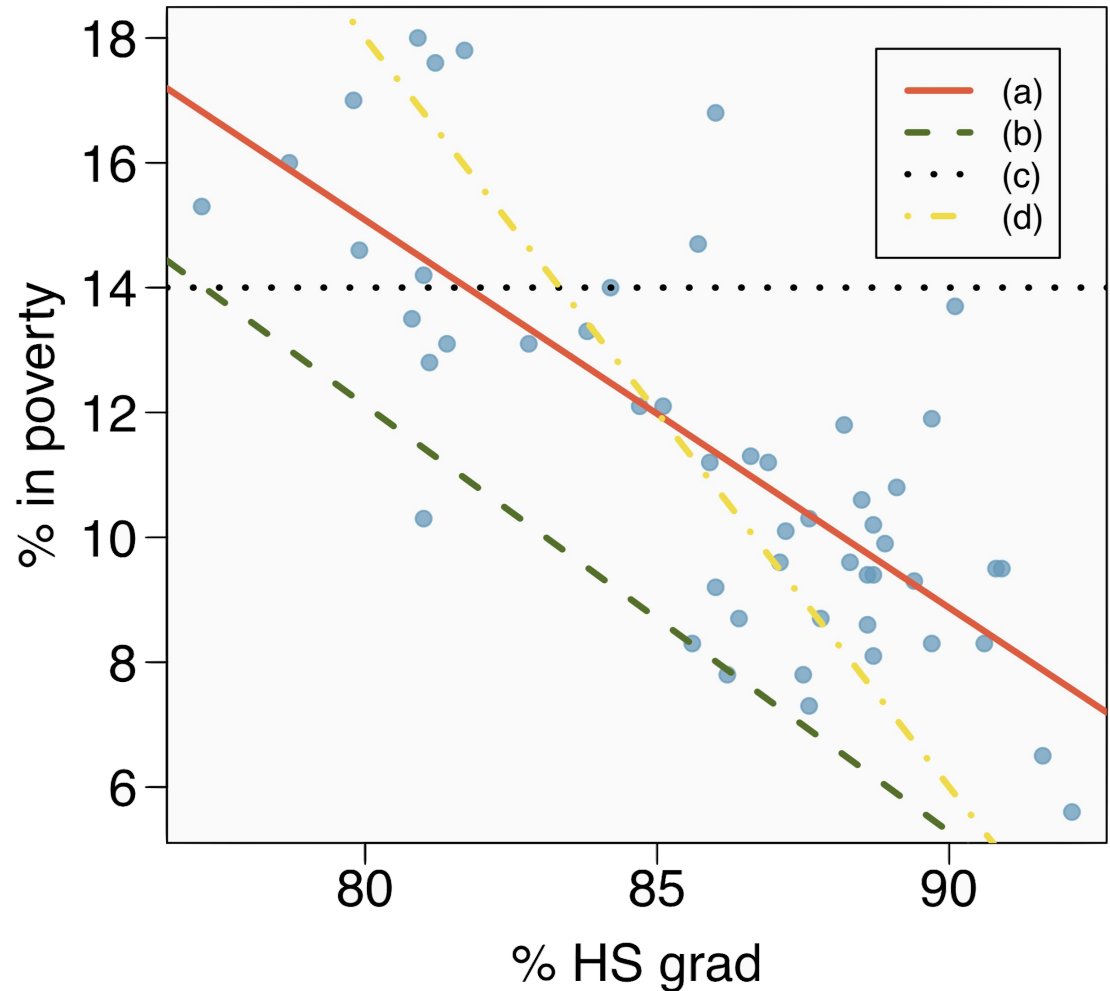
The "hat" is used to signify that this is an estimate.

The high school graduate rate in Georgia is 85.1%. What poverty level does the model predict for this state?

$$64.78 - 0.62 \times 85.1 = 12.018$$

# Eyeballing the line

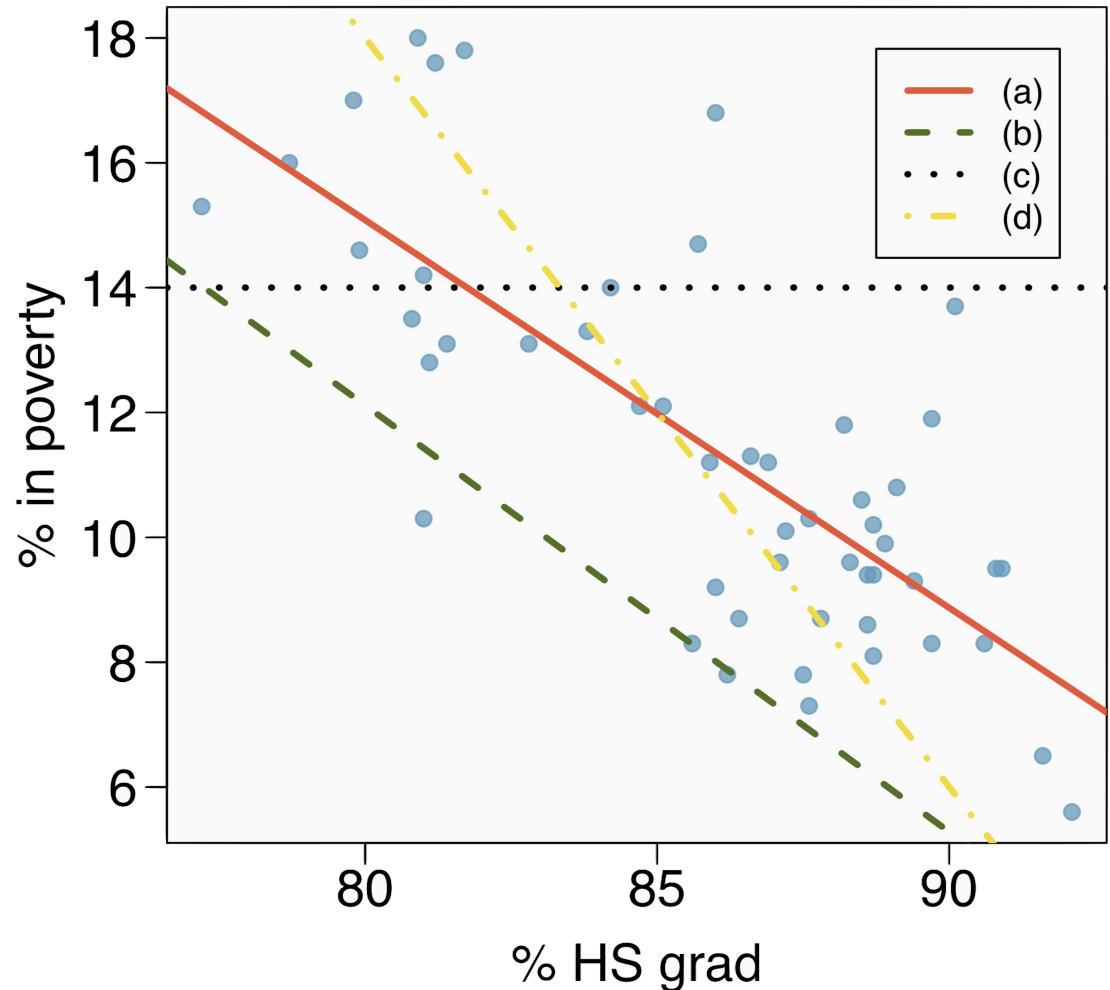
Which of the following appears to be the line that best fits the linear relationship between % in poverty and % HS grad? Choose one.



# Eyeballing the line

Which of the following appears to be the line that best fits the linear relationship between % in poverty and % HS grad? Choose one.

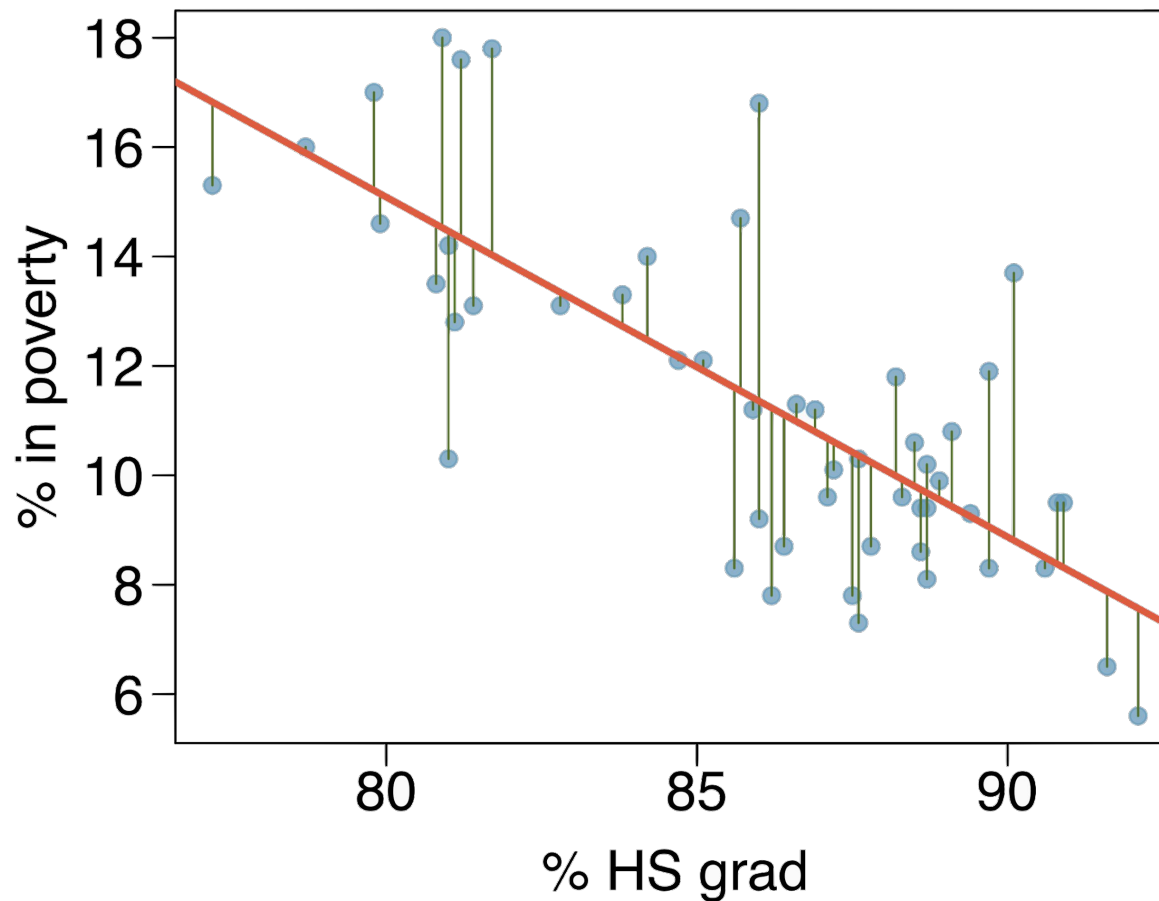
(a)



# Residuals

**Residuals** are the leftovers from the model fit:

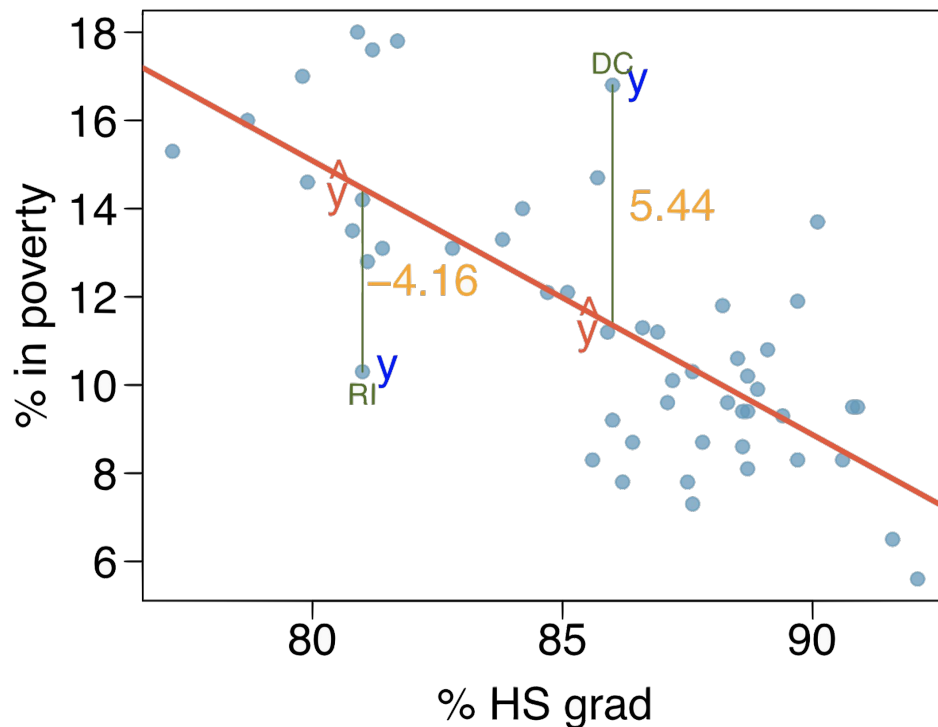
$$\text{Data} = \text{Fit} + \text{Residual}$$



# Residuals (cont.)

Residual is the difference between the observed ( $y_i$ ) and predicted  $\hat{y}_i$ .

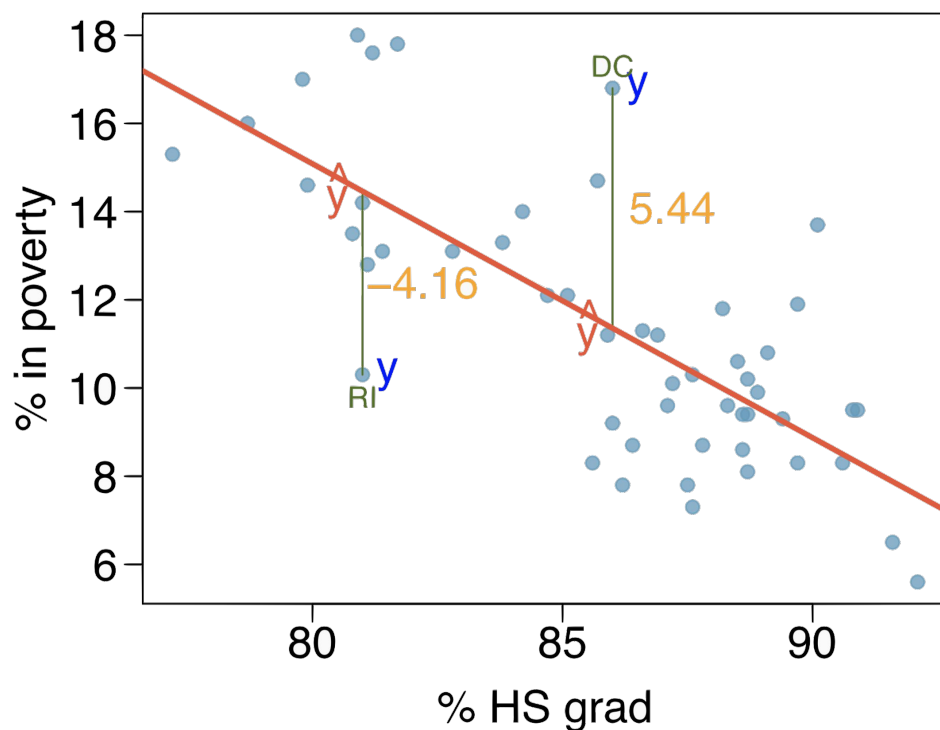
$$e_i = y_i - \hat{y}_i$$



# Residuals (cont.)

Residual is the difference between the observed ( $y_i$ ) and predicted  $\hat{y}_i$ .

$$e_i = y_i - \hat{y}_i$$



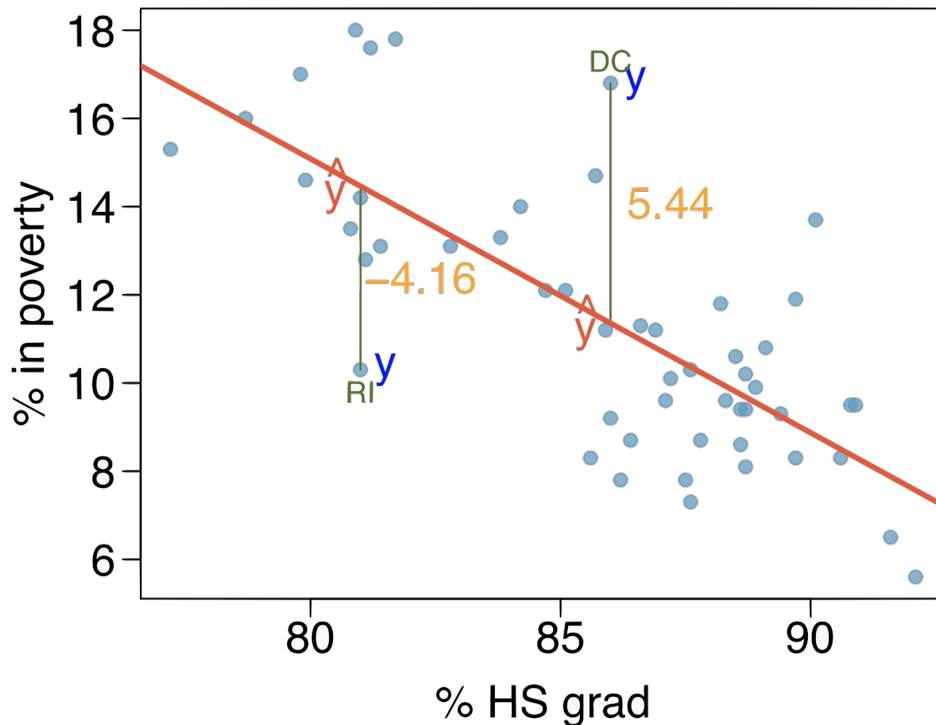
% living in poverty in DC is 5.44% more than predicted.



# Residuals (cont.)

Residual is the difference between the observed ( $y_i$ ) and predicted  $\hat{y}_i$ .

$$e_i = y_i - \hat{y}_i$$



% living in poverty in DC is 5.44% more than predicted.

% living in poverty in RI is 4.16% less than predicted.

# Quantifying the relationship

- *Correlation* describes the strength of the *linear* association between two variables.

# Quantifying the relationship

- *Correlation* describes the strength of the *linear* association between two variables.
- It takes values between -1 (perfect negative) and +1 (perfect positive).

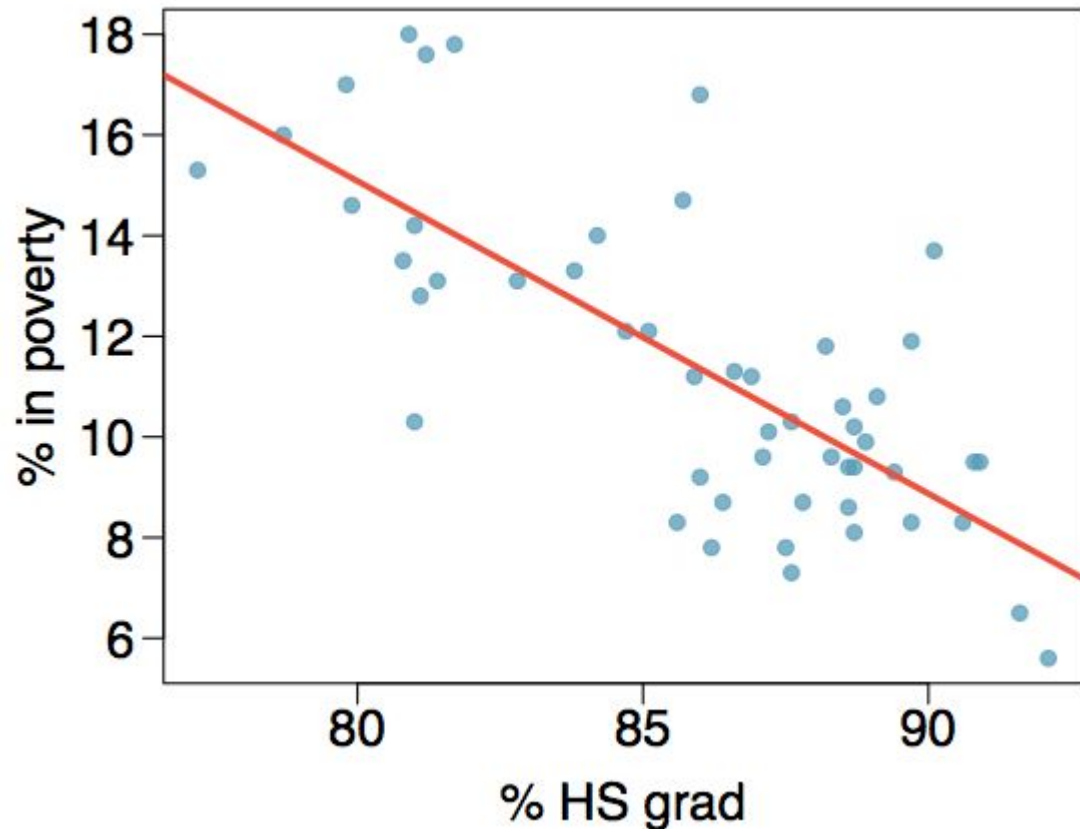
# Quantifying the relationship

- *Correlation* describes the strength of the *linear* association between two variables.
- It takes values between -1 (perfect negative) and +1 (perfect positive).
- A value of 0 indicates no linear association.

# Guessing the correlation

Which of the following is the best guess for the correlation between percent in poverty and percent HS grad?

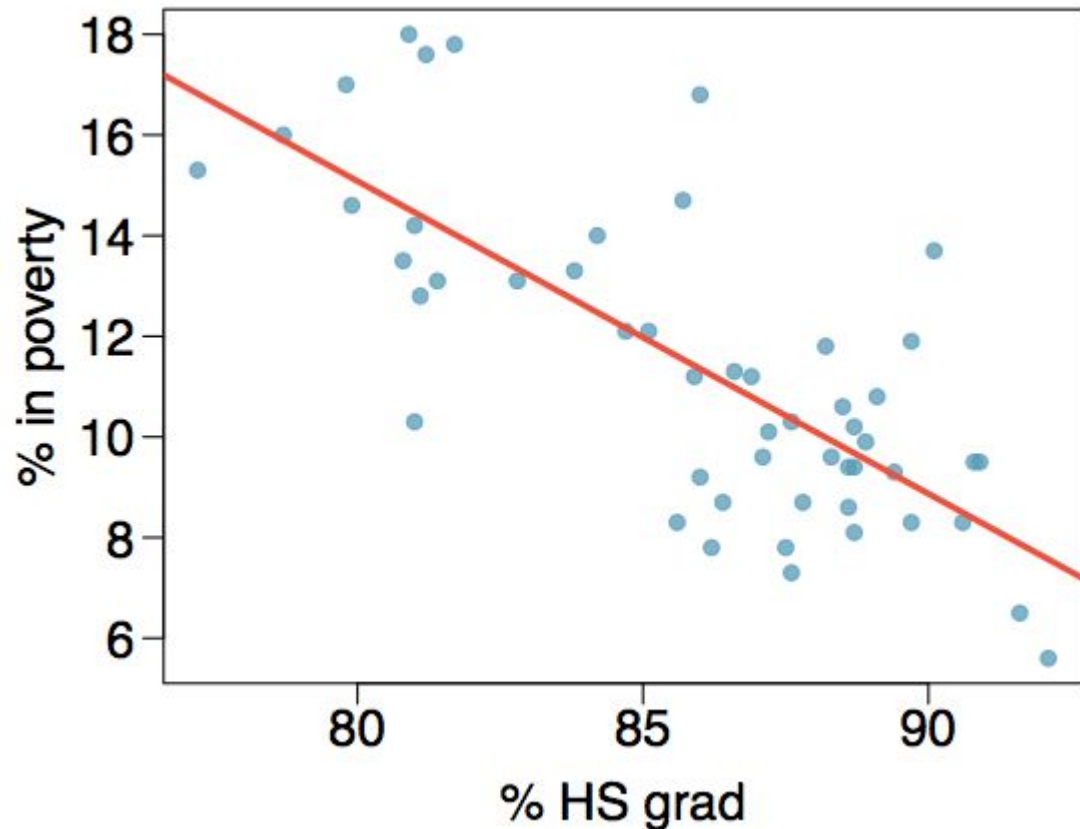
- (a) 0.6
- (b) -0.75
- (c) -0.1
- (d) 0.02
- (e) -1.5



# Guessing the correlation

Which of the following is the best guess for the correlation between percent in poverty and percent HS grad?

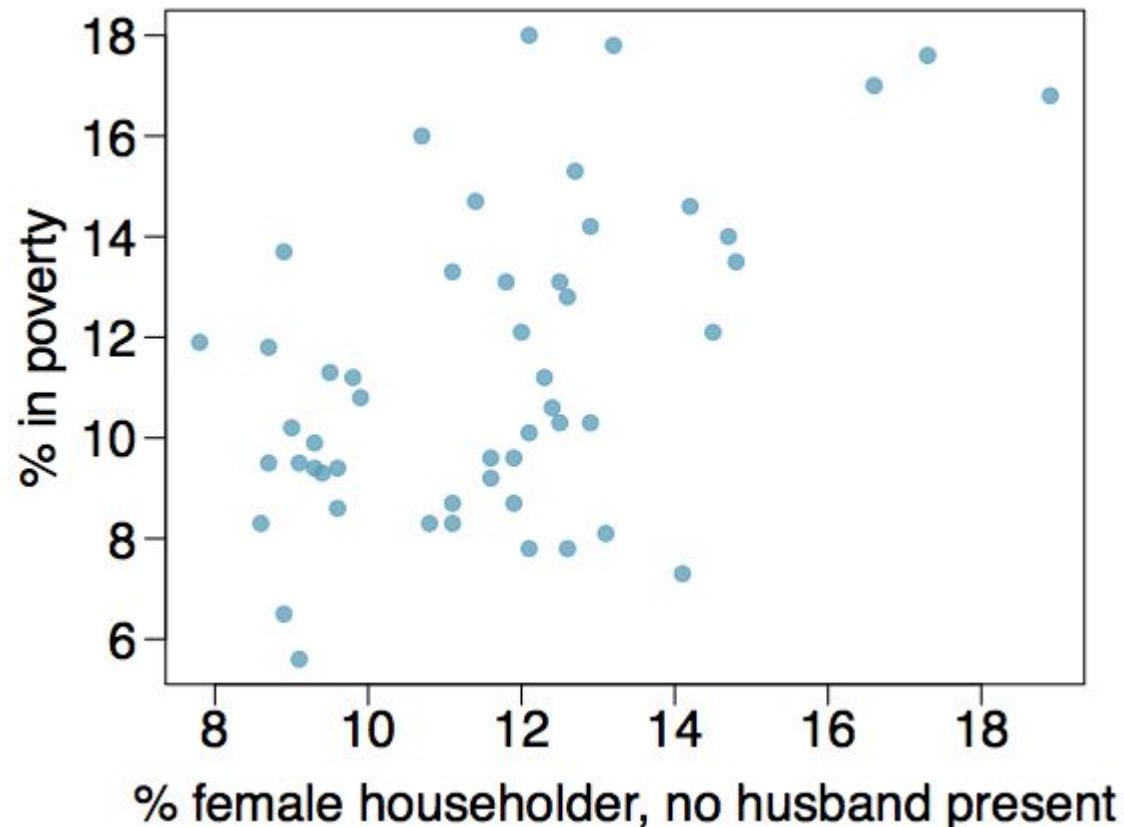
- (a) 0.6
- (b) -0.75*
- (c) -0.1
- (d) 0.02
- (e) -1.5



# Guessing the correlation

Which of the following is the best guess for the correlation between percent in poverty and percent female householder?

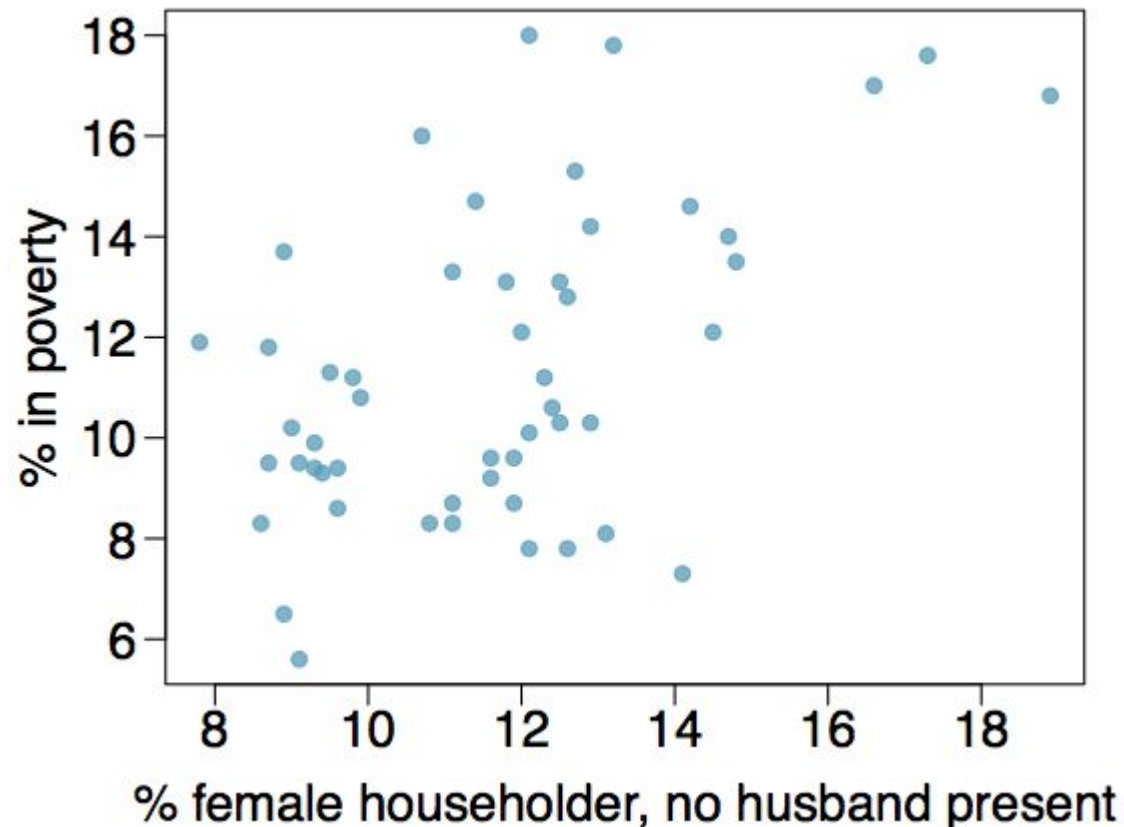
- (a) 0.1
- (b) -0.6
- (c) -0.4
- (d) 0.9
- (e) 0.5



# Guessing the correlation

Which of the following is the best guess for the correlation between percent in poverty and percent female householder?

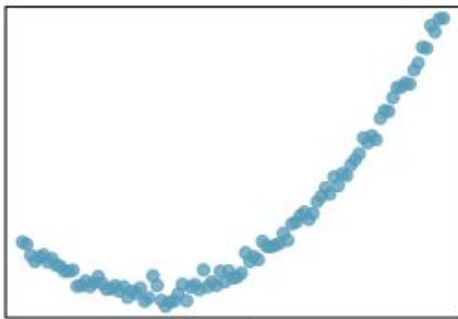
- (a) 0.1
- (b) -0.6
- (c) -0.4
- (d) 0.9
- (e) 0.5



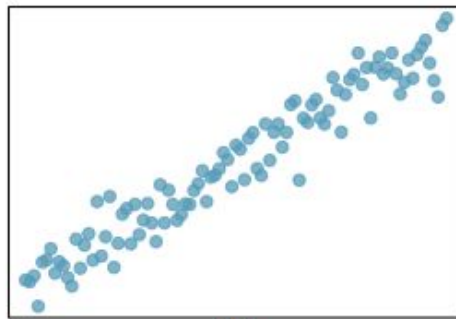


# Assessing the correlation

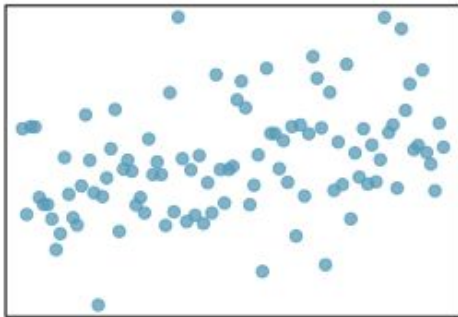
Which of the following has the strongest correlation, i.e. correlation coefficient closest to +1 or -1?



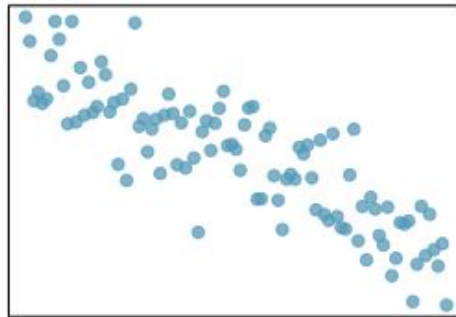
(a)



(b)



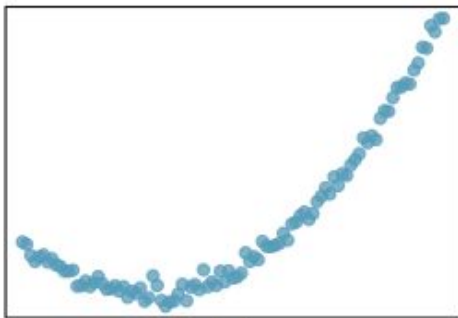
(c)



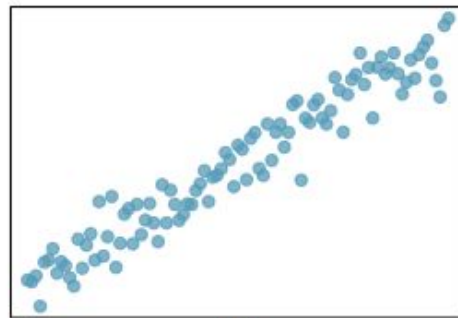
(d)

# Assessing the correlation

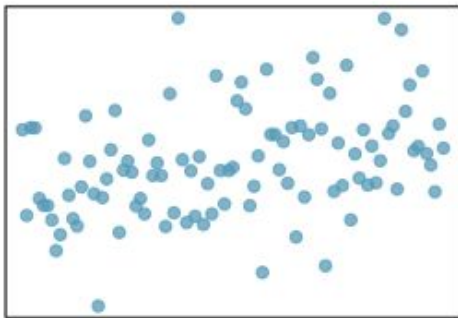
Which of the following has the strongest correlation, i.e. correlation coefficient closest to +1 or -1?



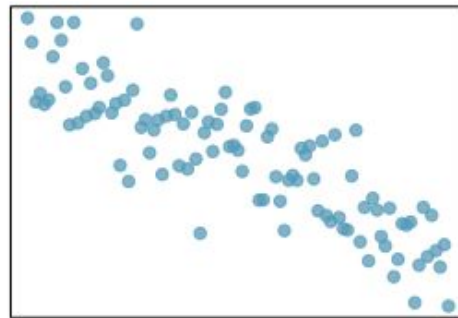
(a)



(b)



(c)



(d)

*(b) → correlation means linear association*

Find more resources at [openintro.org/os](https://openintro.org/os), including

- Slides
- Videos
- Statistical Software Labs
- Discussion Forums (free support for students and teachers)
- Learning Objectives

Teachers only content is also available for [Verified Teachers](#), including

- Exercise solutions
- Sample exams
- Ability to request a free desk copy for a course
- Statistics Teachers email group

Questions? [Contact us](#).