

CAAP Statistics - Lec17

Aug 2, 2022

Review

- Normality Condition
 - Large sample: by CLT
 - Population distribution is normal
 - T-distribution
- One sample mean/proportion
- Paired data
- Difference in two means/proportions
- Equivalence of Hypothesis Testing and Confidence Interval

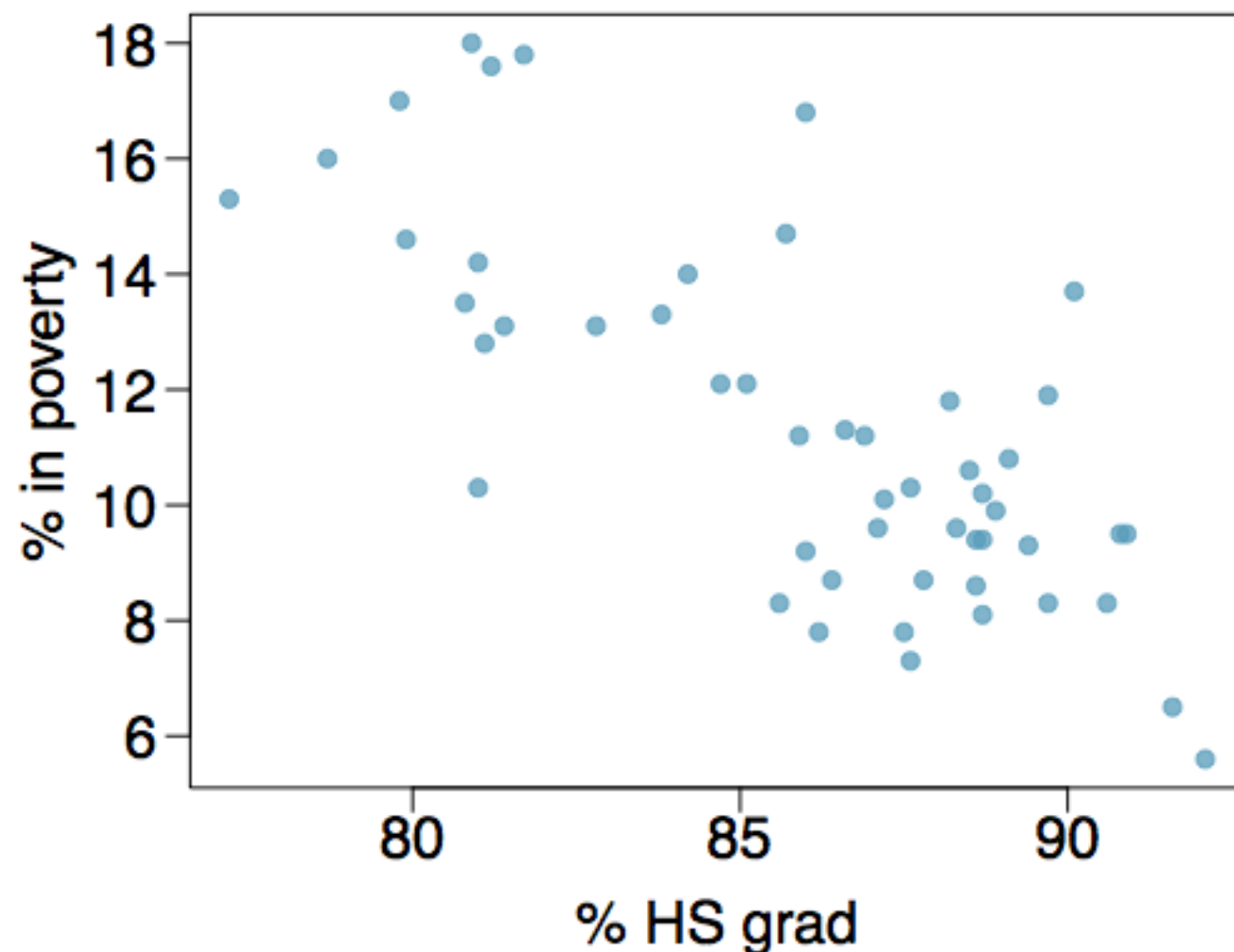
Learning Objectives

- Line Fitting, Residuals and Correlation
- Fitting a line by Least Squares Regression
- Types of outliers in Linear Regression
 - Outliers, High leverage point
 - Influential point
- Inference for Linear Regression
 - Hypothesis testing for the slope

Line Fitting, Residuals, and Correlation

Poverty vs. HS graduate rate

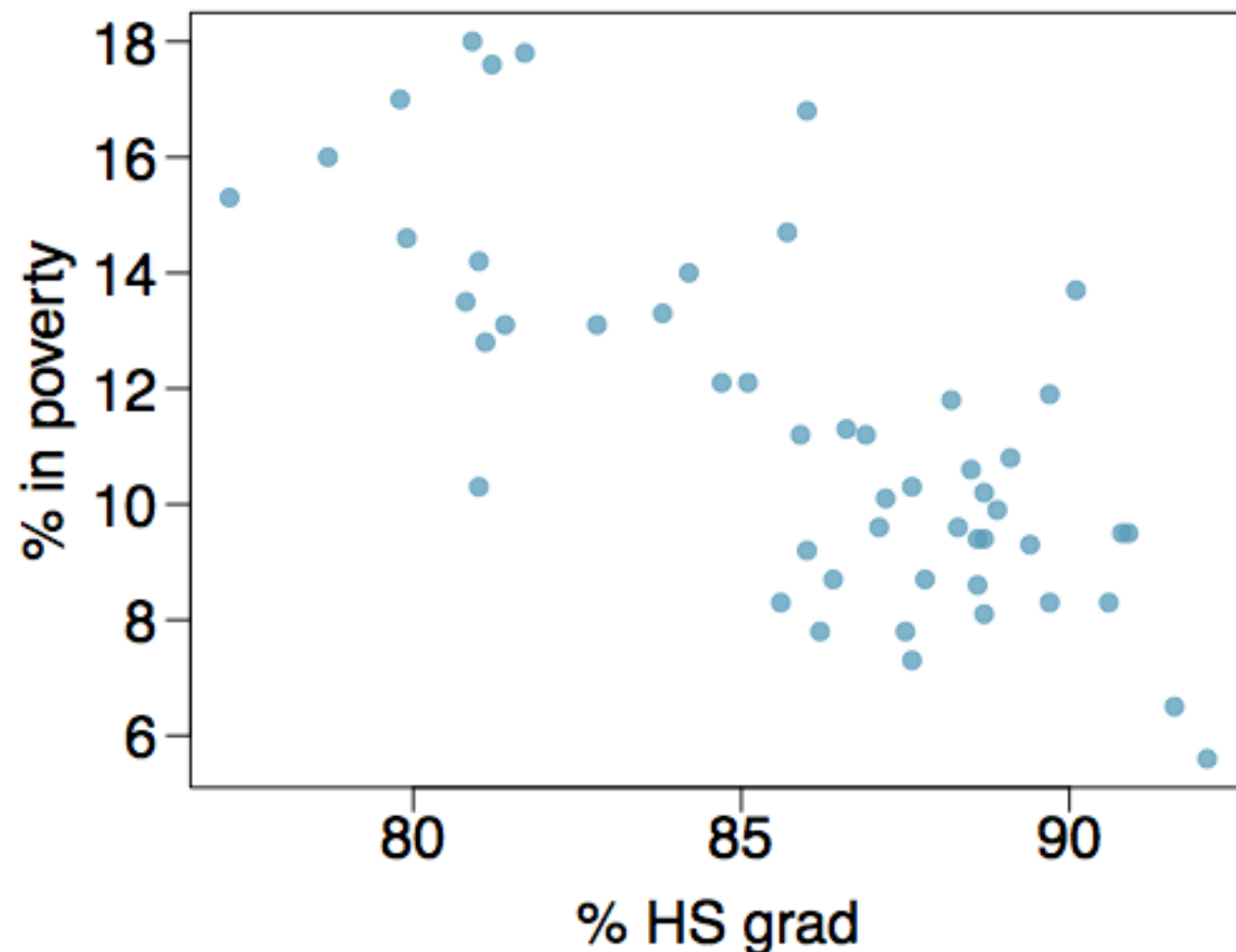
The *scatterplot* below shows the relationship between HS graduate rate in all 50 US states and DC and the percent of residents who live below the poverty line (income below \$23,050 for a family of 4 in 2012).



Response variable?

Poverty vs. HS graduate rate

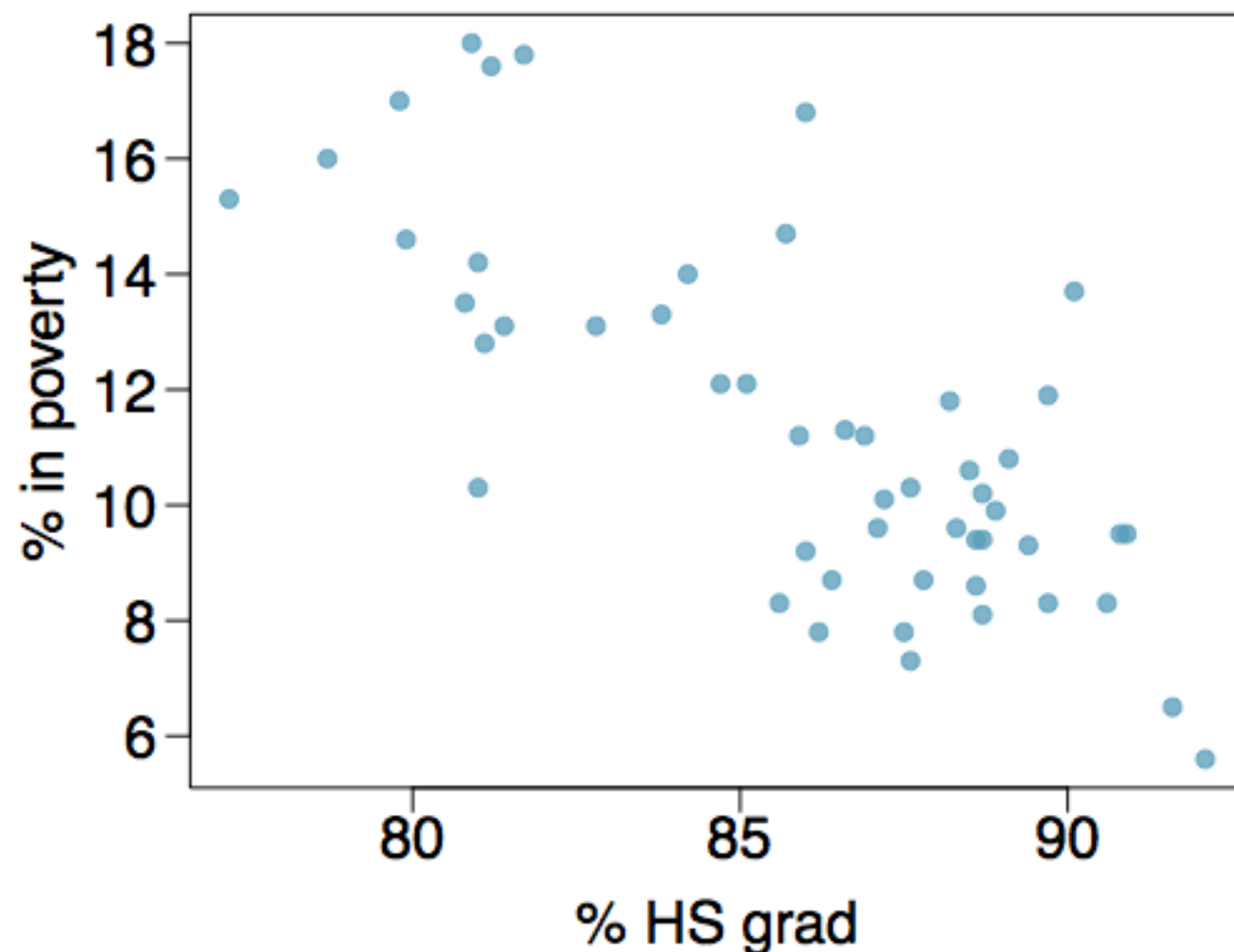
The *scatterplot* below shows the relationship between HS graduate rate in all 50 US states and DC and the percent of residents who live below the poverty line (income below \$23,050 for a family of 4 in 2012).



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).



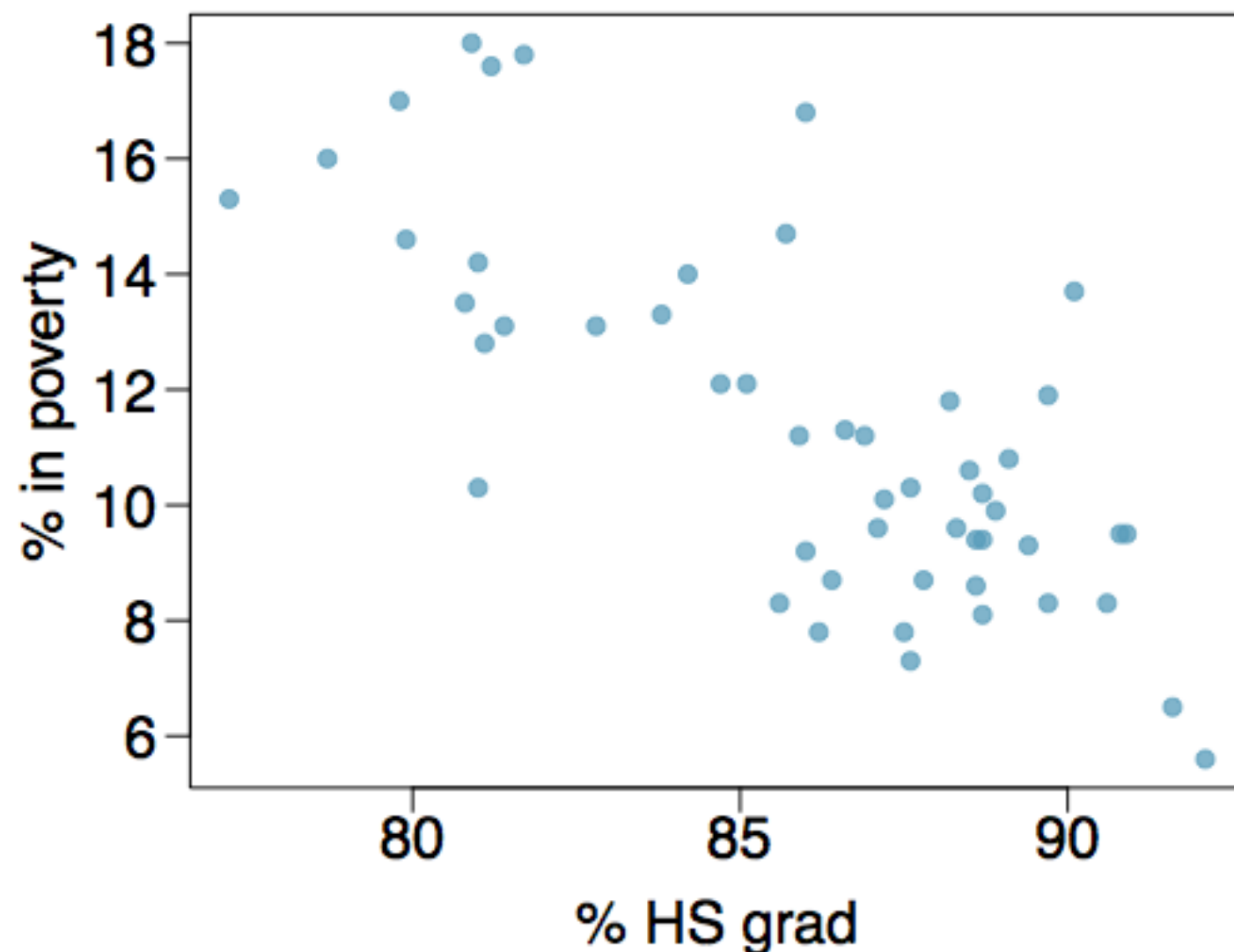
Response variable?

% 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?

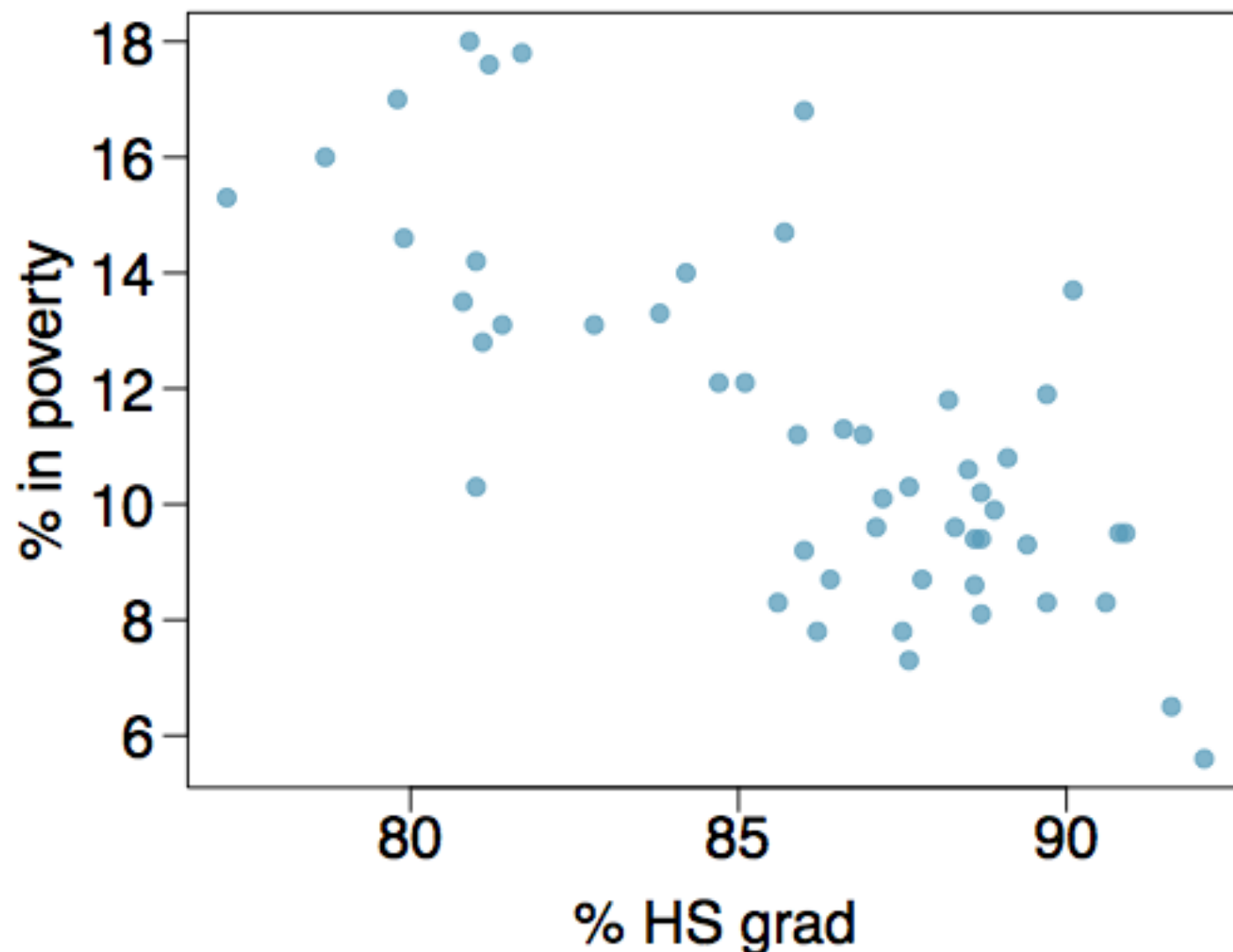
% in poverty

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

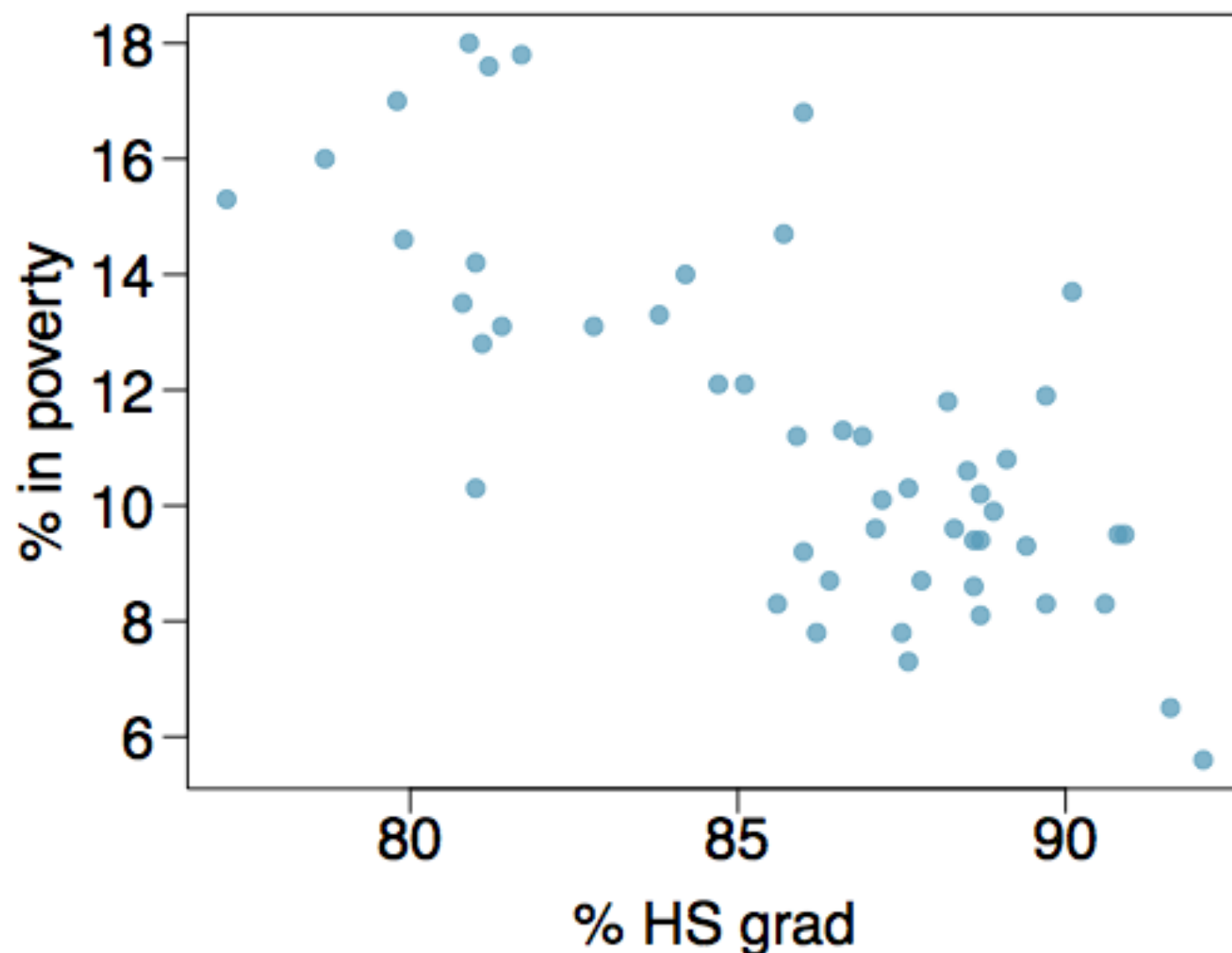
Explanatory variable?

% 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

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

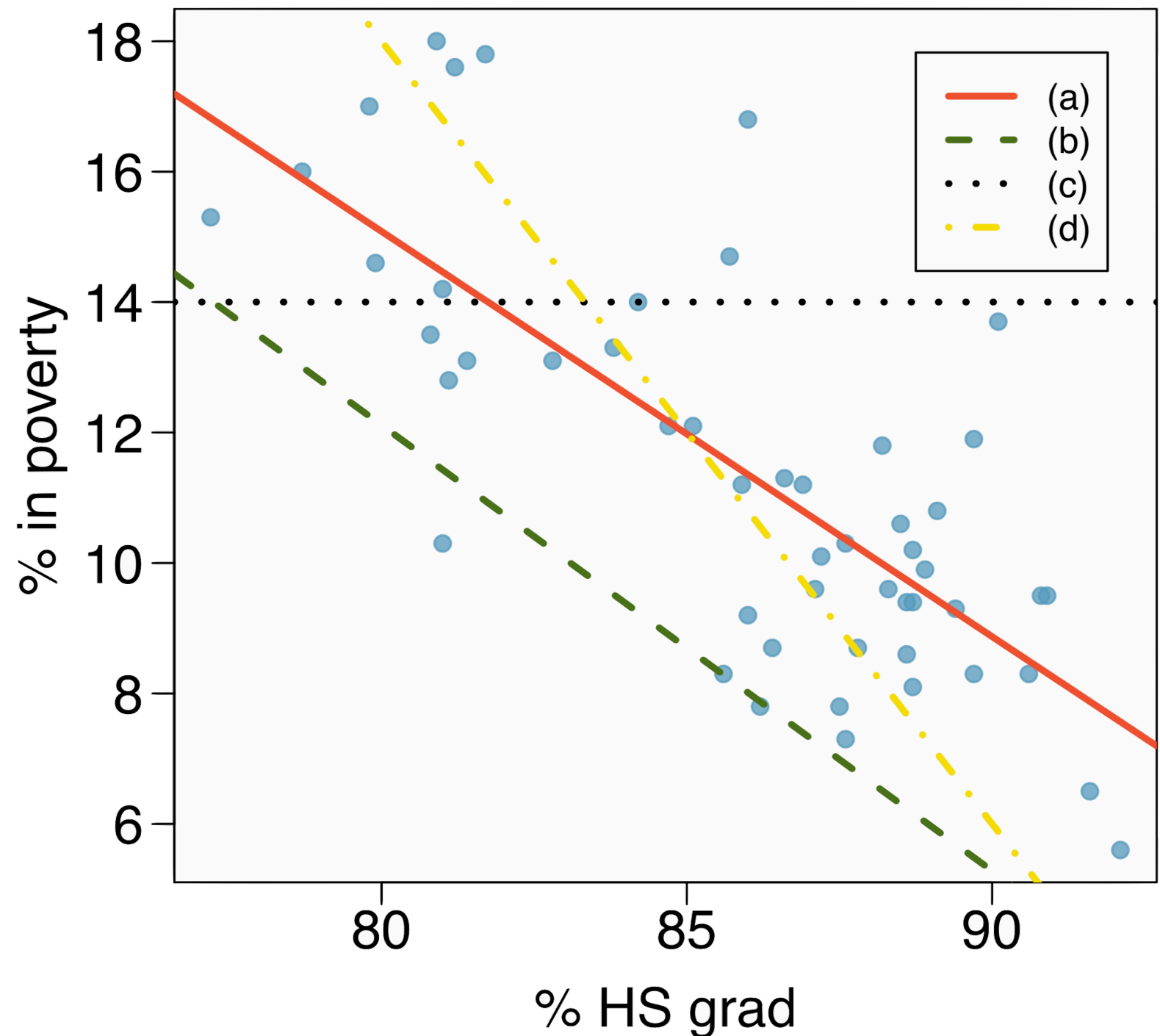
The "hat" is used to signify that this is an estimate.

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

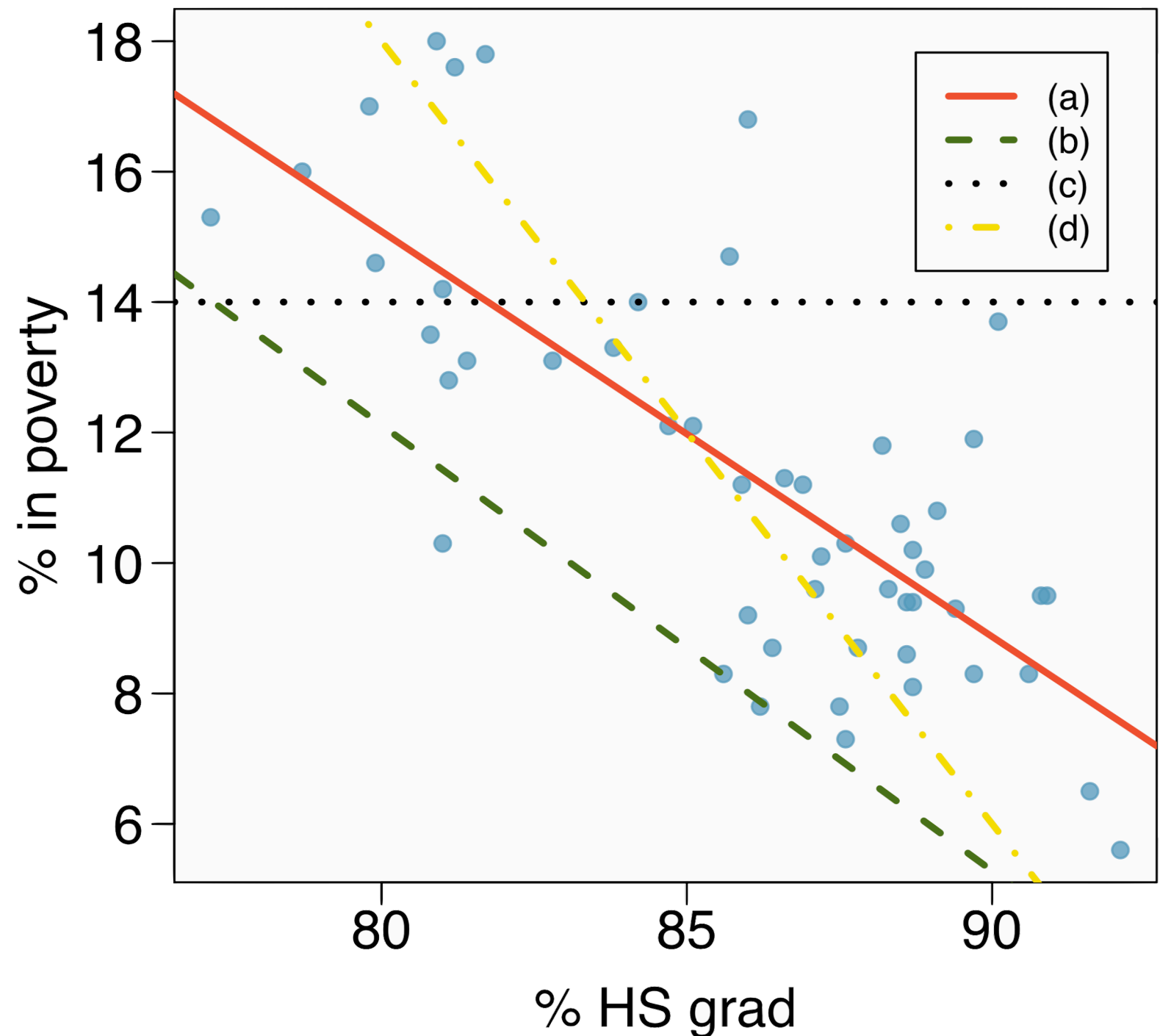
Which of the following appears to be the line that best fits the linear relationship between % in poverty and % HS grad? Choose one.



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.

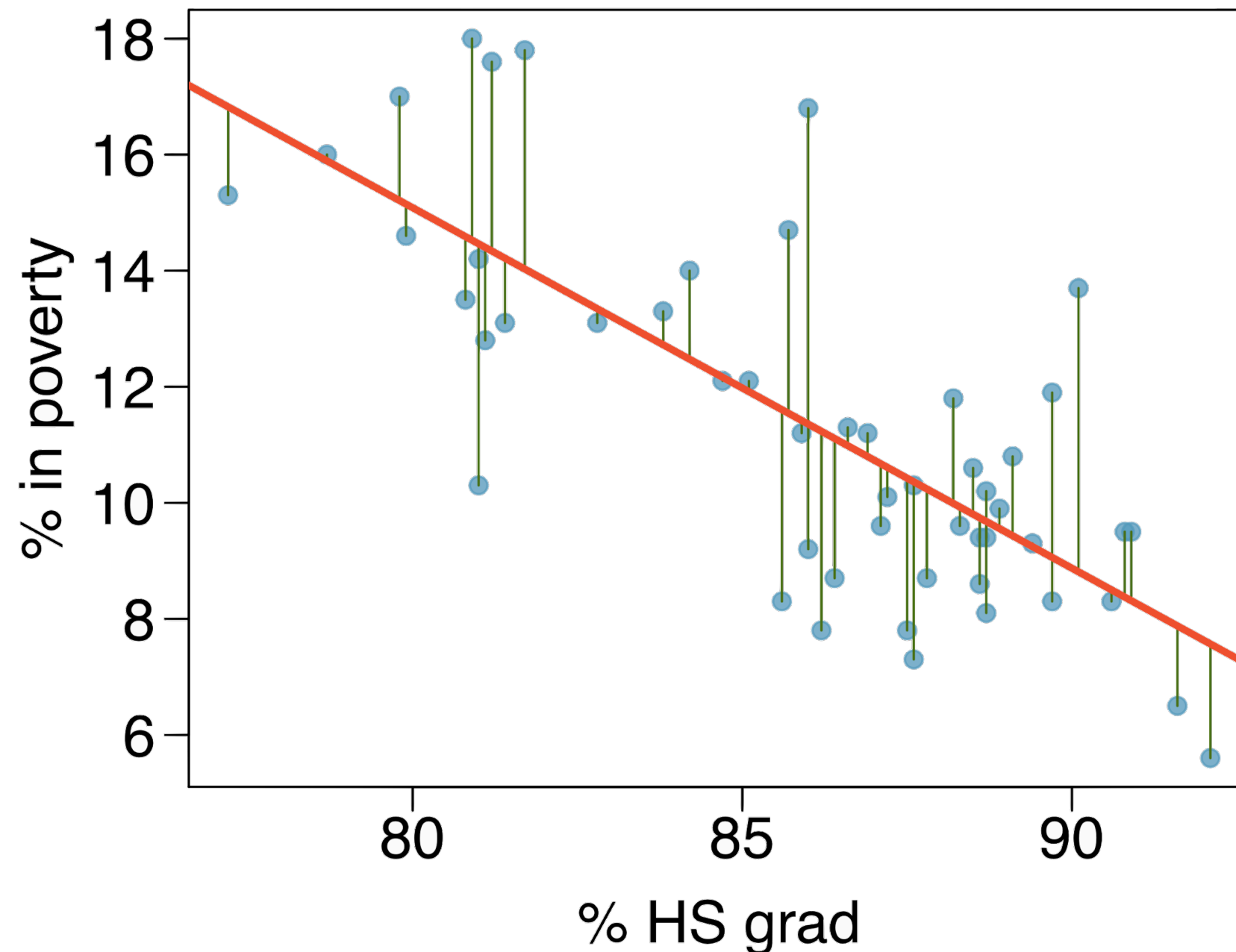
(a)



Residuals

Residuals are the leftovers from the model fit:

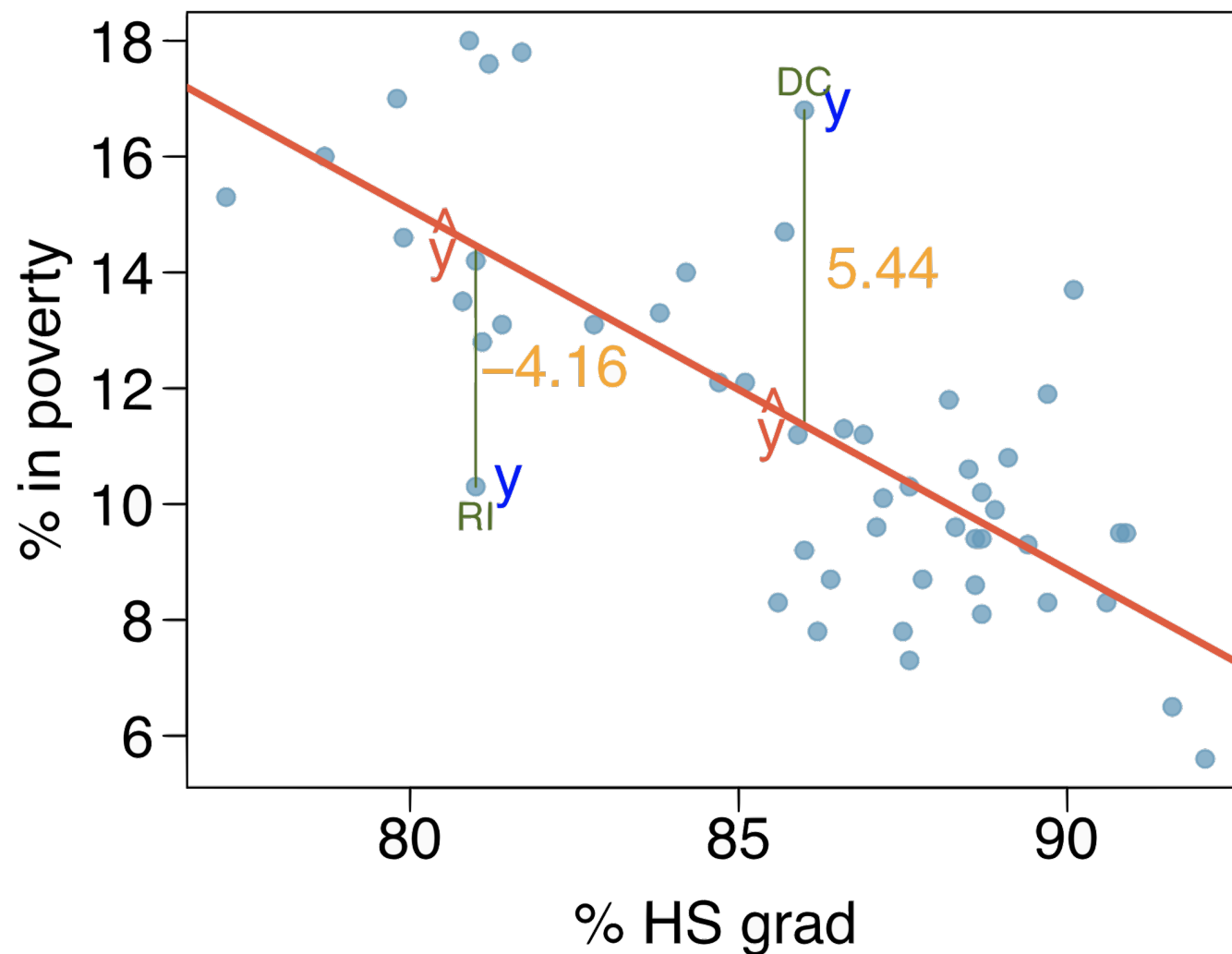
$$\text{Data} = \text{Fit} + \text{Residual}$$



Residuals (cont.)

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

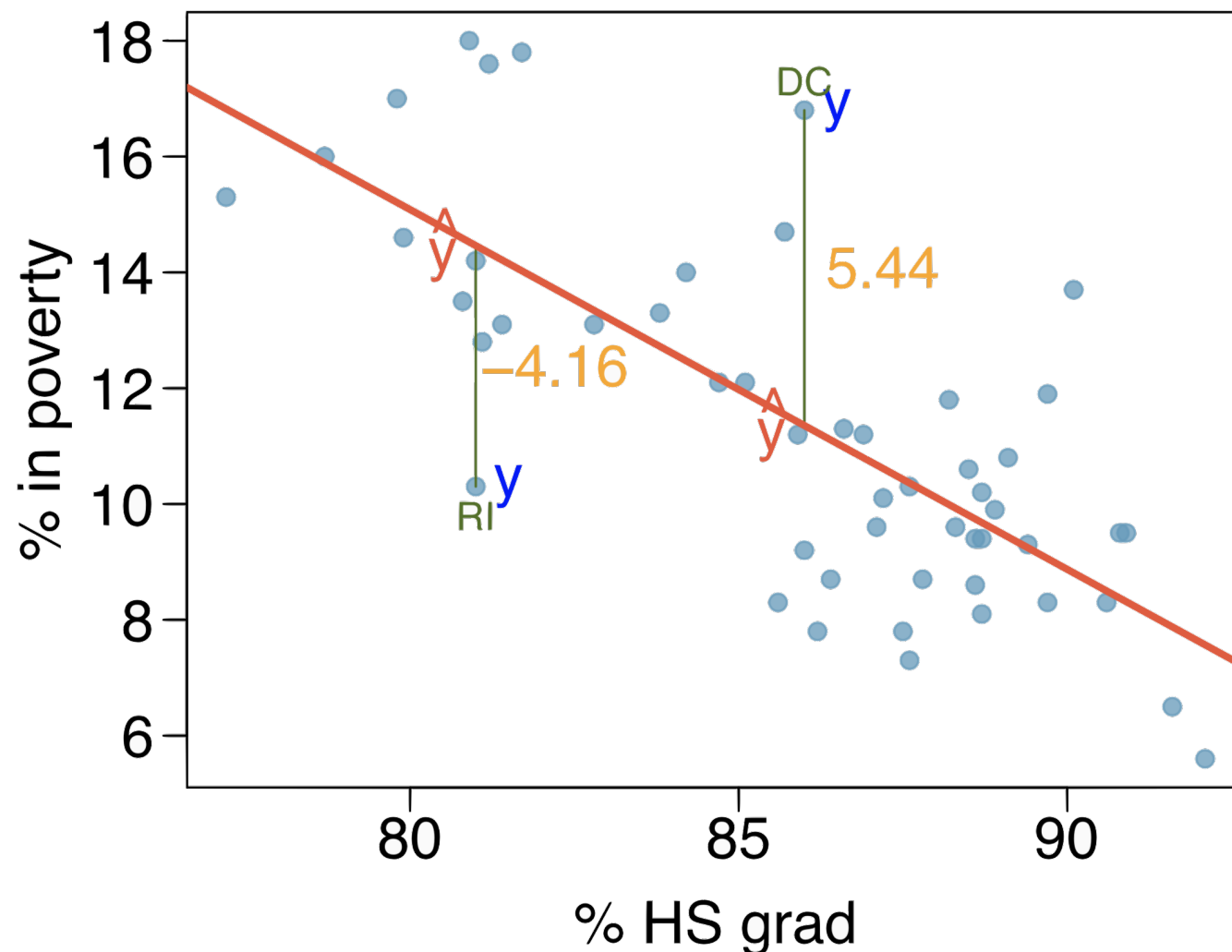
$$e_i = y_i - \hat{y}_i$$



Residuals (cont.)

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

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

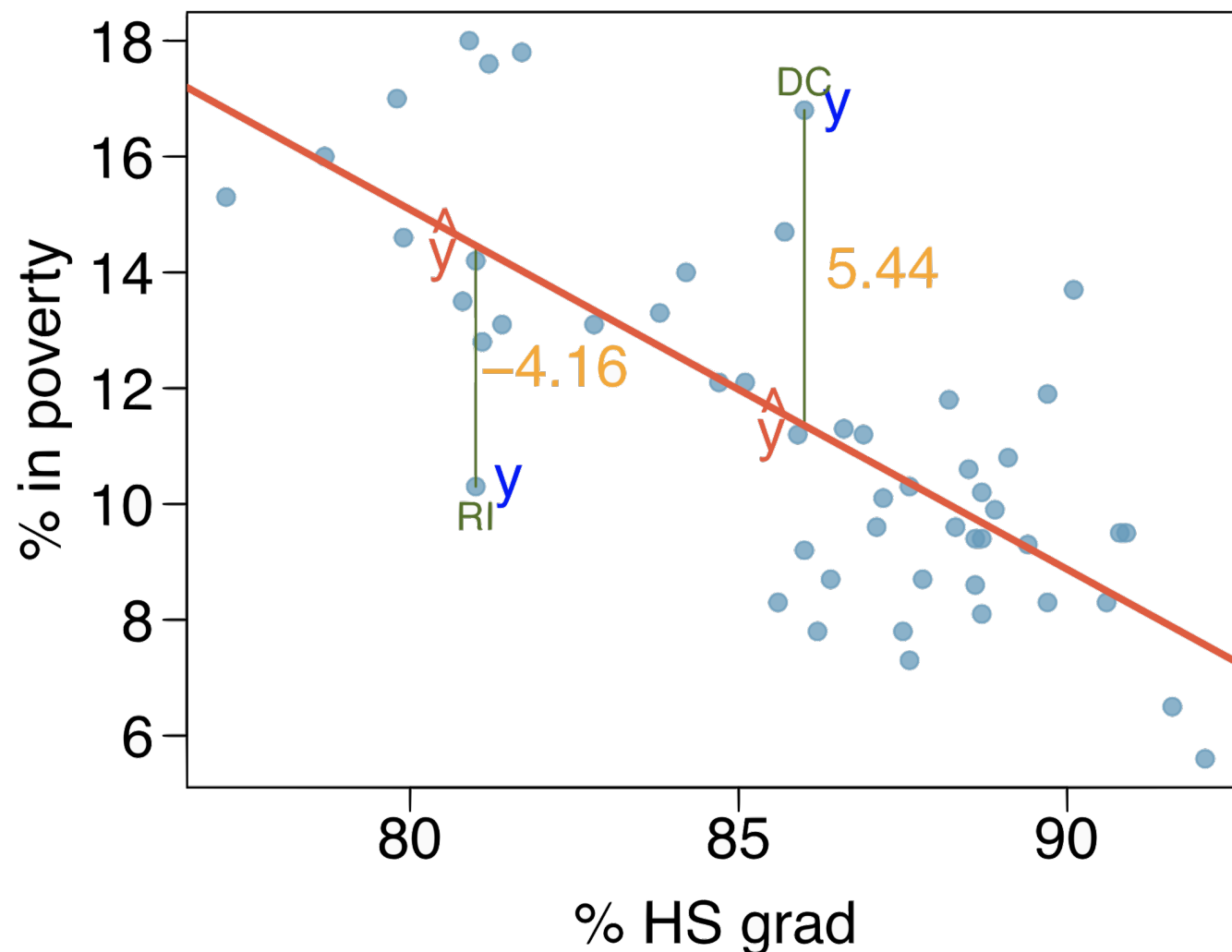


% living in poverty in DC is 5.44% more than predicted.

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).

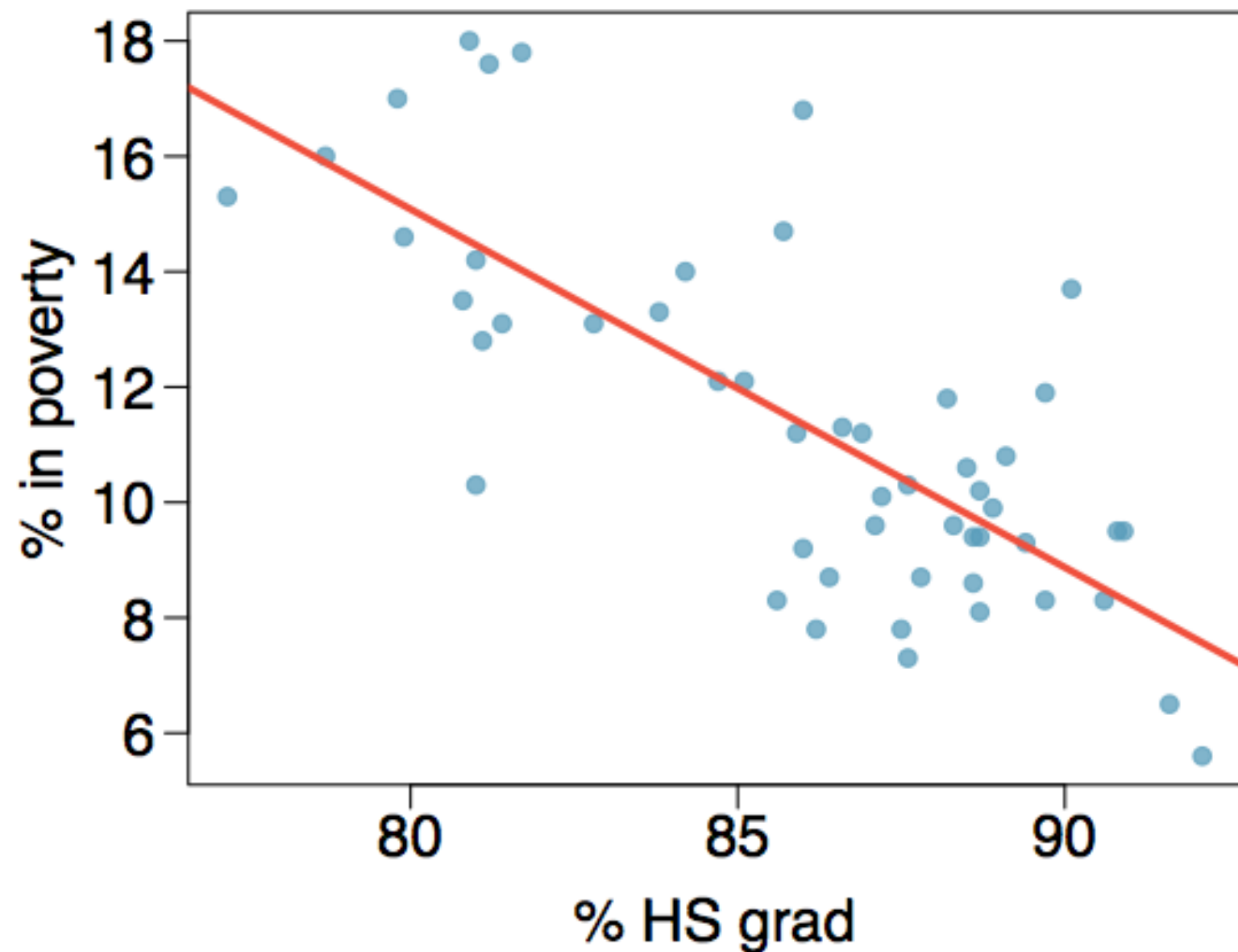
Quantifying the relationship

- *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?

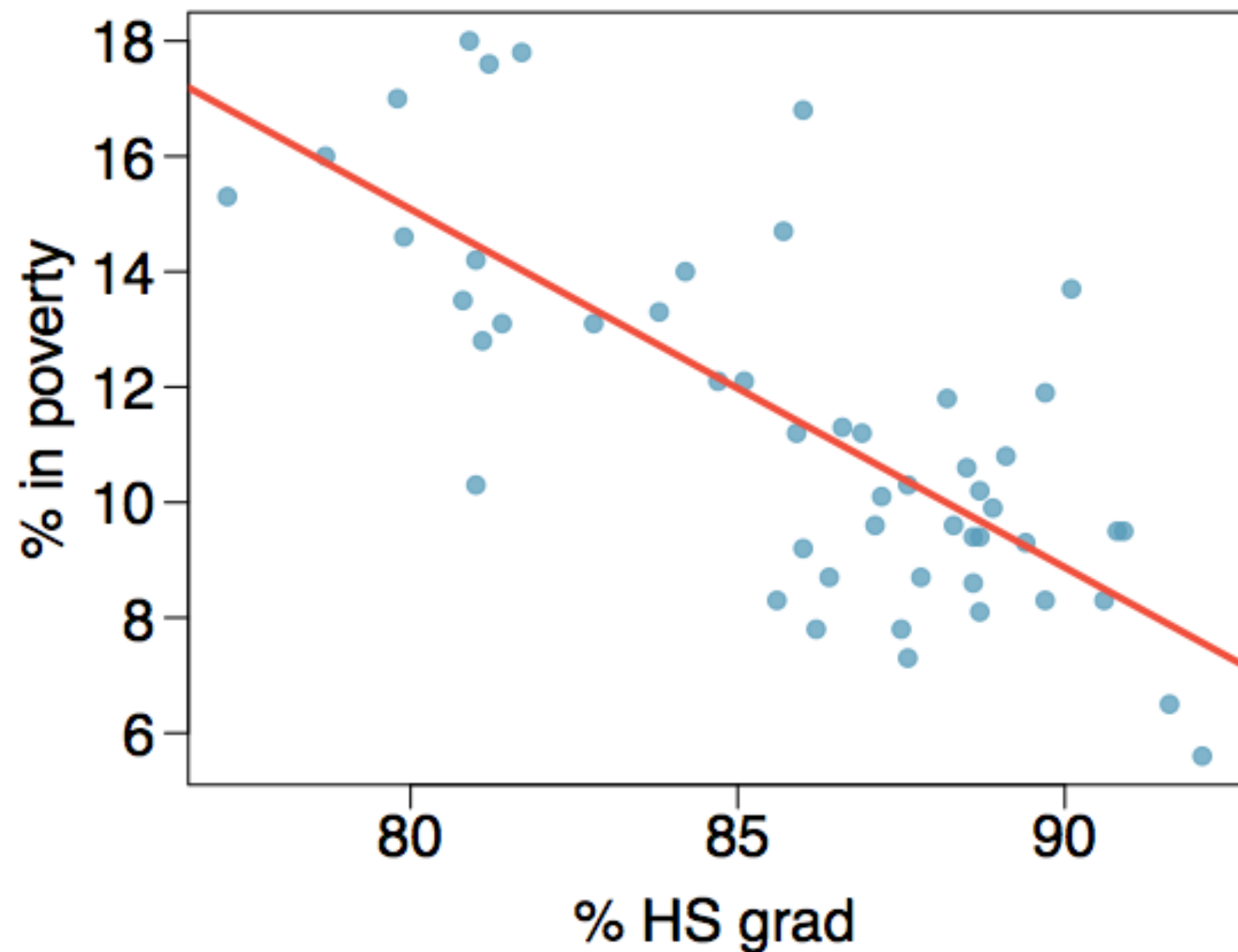
- (a) 0.6
- (b) -0.75
- (c) -0.1
- (d) 0.02
- (e) -1.5



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

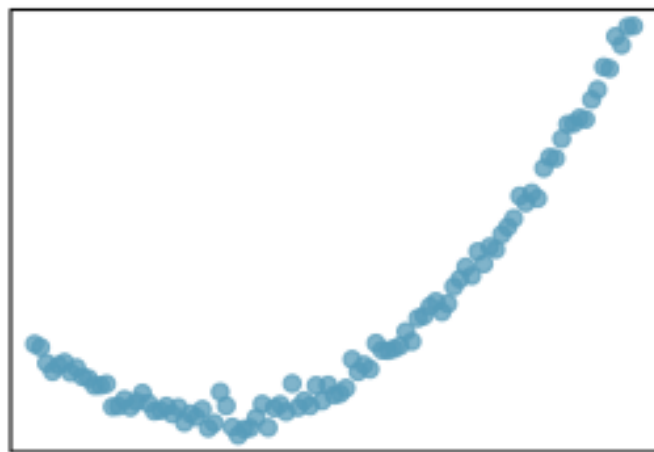
Which of the following is the best guess for the correlation between percent in poverty and percent female householder?

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

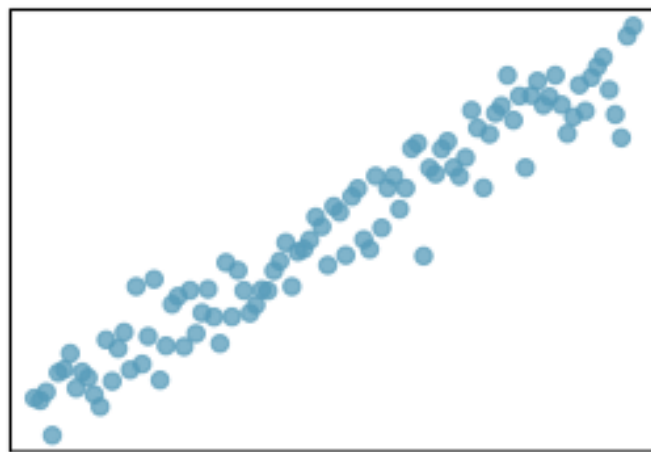


Assessing the correlation

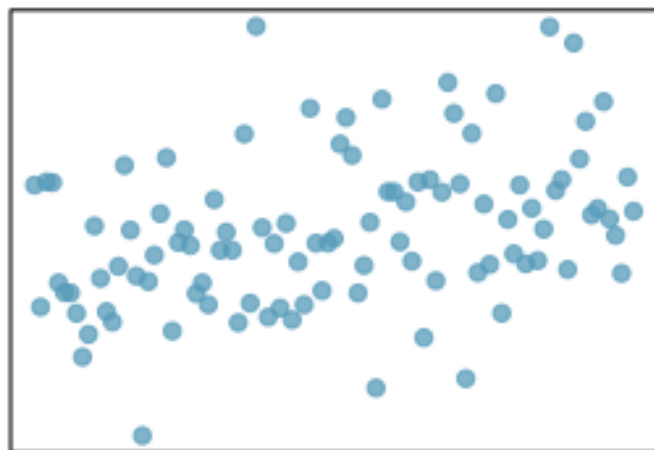
Which of the following is has the strongest correlation, i.e. correlation coefficient closest to +1 or -1?



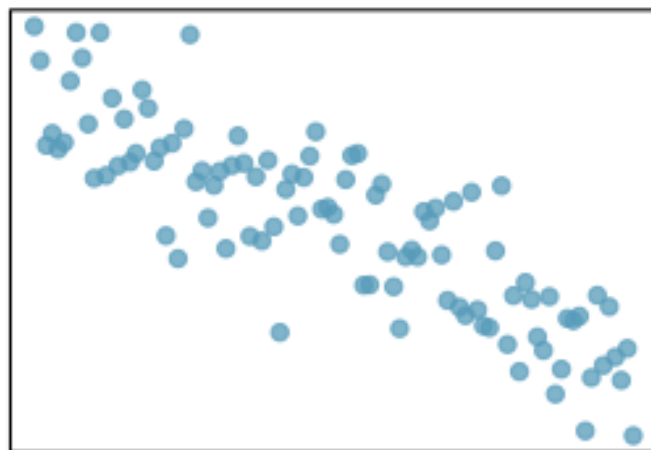
(a)



(b)



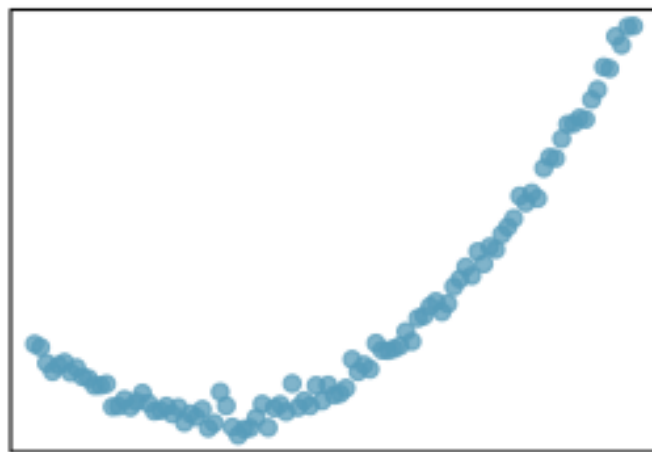
(c)



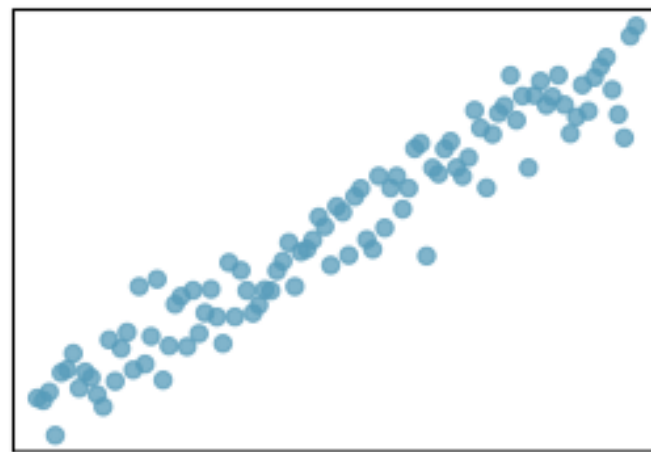
(d)

Assessing the correlation

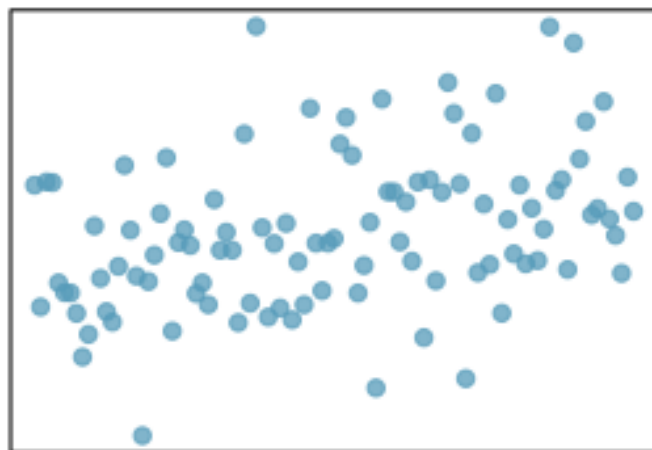
Which of the following is has the strongest correlation, i.e. correlation coefficient closest to +1 or -1?



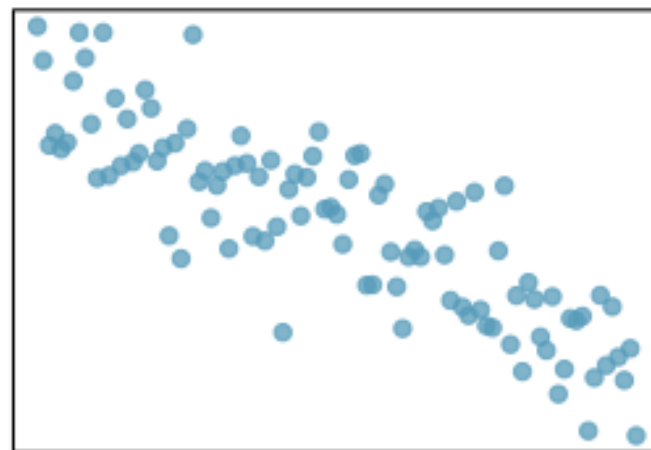
(a)



(b)



(c)



(d)

*(b) → correlation means
linear association*

Fitting a line by least squares regression

A measure for the best line

- We want a line that has small residuals

A measure for the best line

- We want a line that has small residuals
 1. Option 1: Minimize the sum of magnitudes (absolute values) of residuals

$$|e_1| + |e_2| + \dots + |e_n|$$

A measure for the best line

- We want a line that has small residuals
 1. Option 1: Minimize the sum of magnitudes (absolute values) of residuals

$$|e_1| + |e_2| + \dots + |e_n|$$

2. Option 2: Minimize the sum of squared residuals -- *least squares*

$$e_1^2 + e_2^2 + \dots + e_n^2$$

A measure for the best line

- We want a line that has small residuals
 1. Option 1: Minimize the sum of magnitudes (absolute values) of residuals

$$|e_1| + |e_2| + \dots + |e_n|$$

2. Option 2: Minimize the sum of squared residuals -- *least squares*

$$e_1^2 + e_2^2 + \dots + e_n^2$$

- Why least squares?

A measure for the best line

- We want a line that has small residuals
 1. Option 1: Minimize the sum of magnitudes (absolute values) of residuals

$$|e_1| + |e_2| + \dots + |e_n|$$

2. Option 2: Minimize the sum of squared residuals -- *least squares*

$$e_1^2 + e_2^2 + \dots + e_n^2$$

- Why least squares?
 1. Most commonly used

A measure for the best line

- We want a line that has small residuals
 1. Option 1: Minimize the sum of magnitudes (absolute values) of residuals

$$|e_1| + |e_2| + \dots + |e_n|$$

2. Option 2: Minimize the sum of squared residuals -- *least squares*

$$e_1^2 + e_2^2 + \dots + e_n^2$$

- Why least squares?
 1. Most commonly used
 2. Easier to compute by hand and using software

A measure for the best line

- We want a line that has small residuals
 1. Option 1: Minimize the sum of magnitudes (absolute values) of residuals

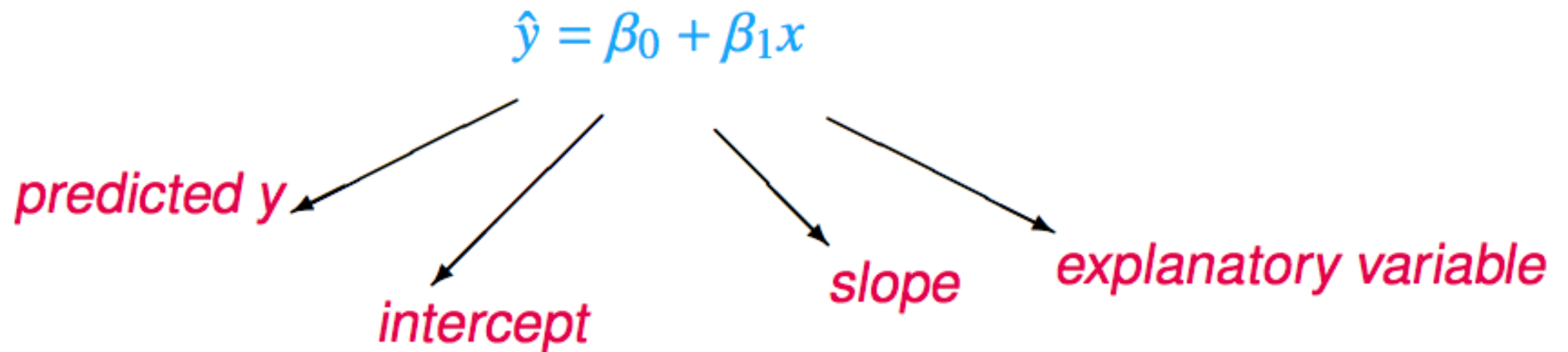
$$|e_1| + |e_2| + \dots + |e_n|$$

2. Option 2: Minimize the sum of squared residuals -- *least squares*

$$e_1^2 + e_2^2 + \dots + e_n^2$$

- Why least squares?
 1. Most commonly used
 2. Easier to compute by hand and using software
 3. In many applications, a residual twice as large as another is usually more than twice as bad

The least squares line



Notation:

- Intercept:
 - Parameter: β_0
 - Point estimate: b_0
- Slope:
 - Parameter: β_1
 - Point estimate: b_1

Conditions for the least squares line

1. Linearity

Conditions for the least squares line

1. Linearity

2. Nearly normal residuals

Conditions for the least squares line

1. Linearity

2. Nearly normal residuals

3. Constant variability

Conditions: (1) Linearity

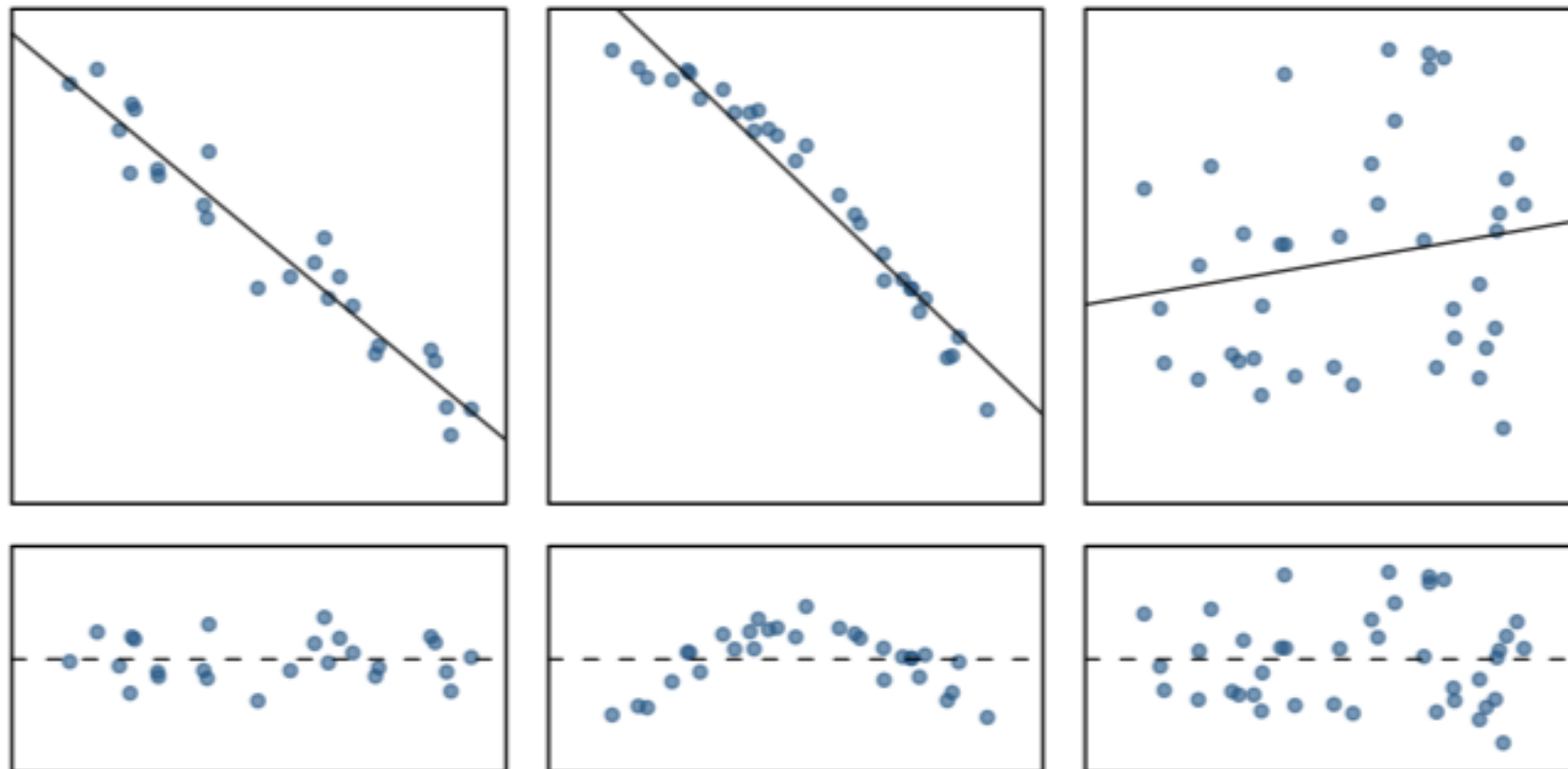
- The relationship between the explanatory and the response variable should be linear.

Conditions: (1) Linearity

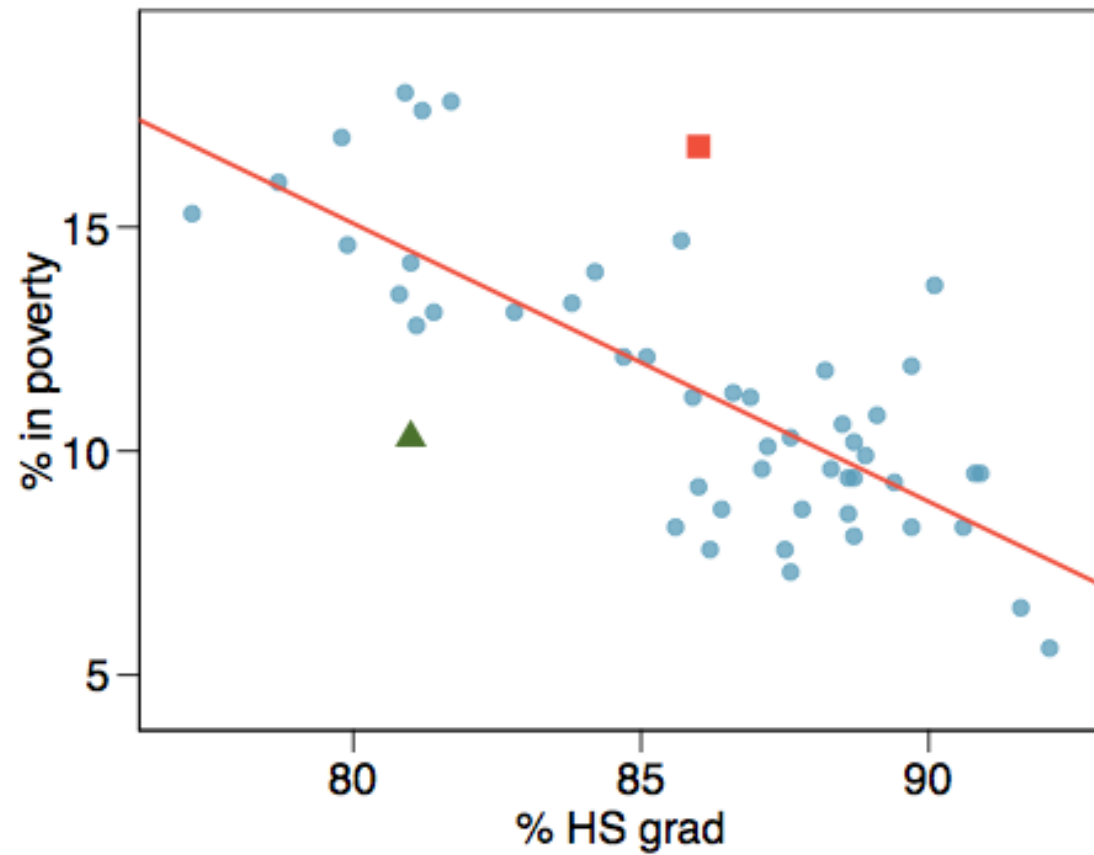
- The relationship between the explanatory and the response variable should be linear.
- Methods for fitting a model to non-linear relationships exist, but are beyond the scope of this class.

Conditions: (1) Linearity

- The relationship between the explanatory and the response variable should be linear.
- Methods for fitting a model to non-linear relationships exist, but are beyond the scope of this class.
- Check using a scatterplot of the data, or a *residuals plot*.

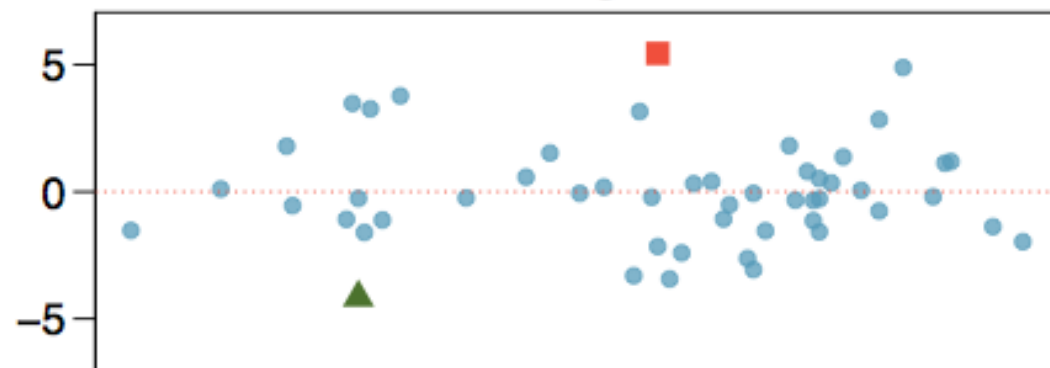


Anatomy of a residuals plot

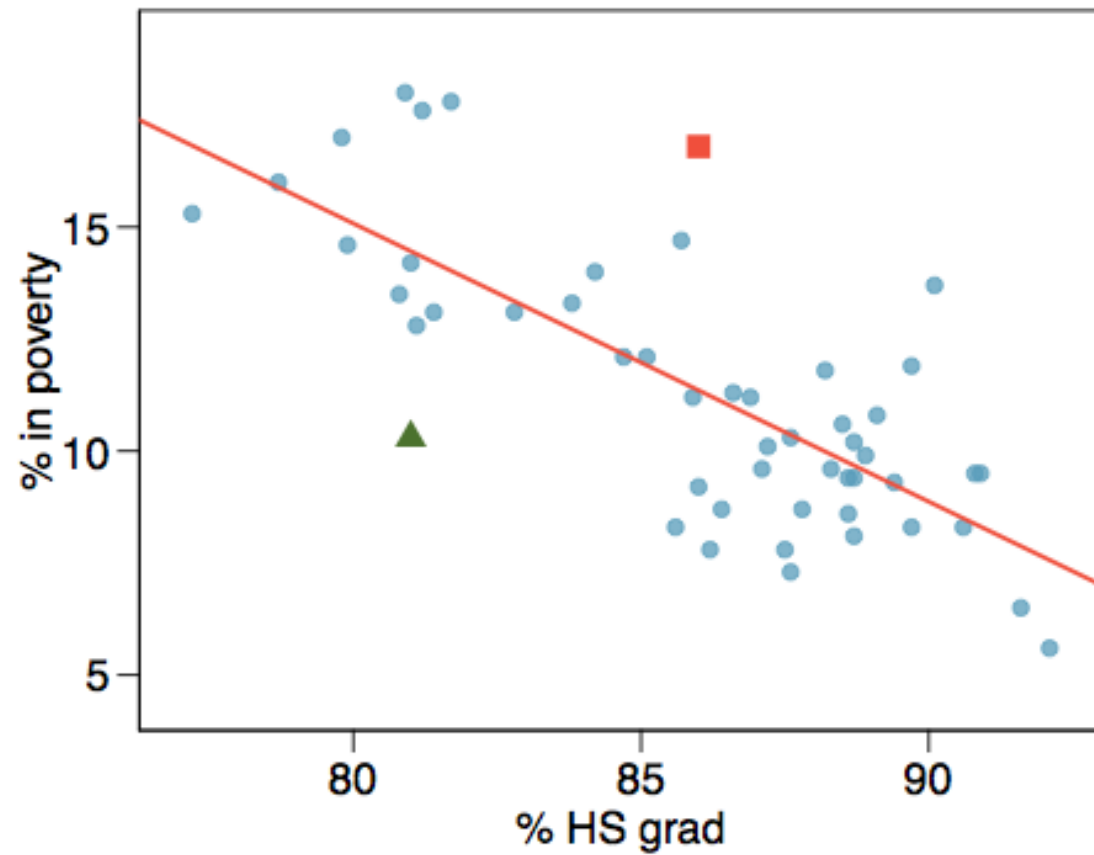


▲ RI:

$$\begin{aligned}\% \text{ HS grad} &= 81 & \% \text{ in poverty} &= 10.3 \\ \% \text{ in } \widehat{\text{poverty}} &= 64.68 - 0.62 * 81 = 14.46 \\ e &= \% \text{ in poverty} - \% \text{ in } \widehat{\text{poverty}} \\ &= 10.3 - 14.46 = -4.16\end{aligned}$$



Anatomy of a residuals plot

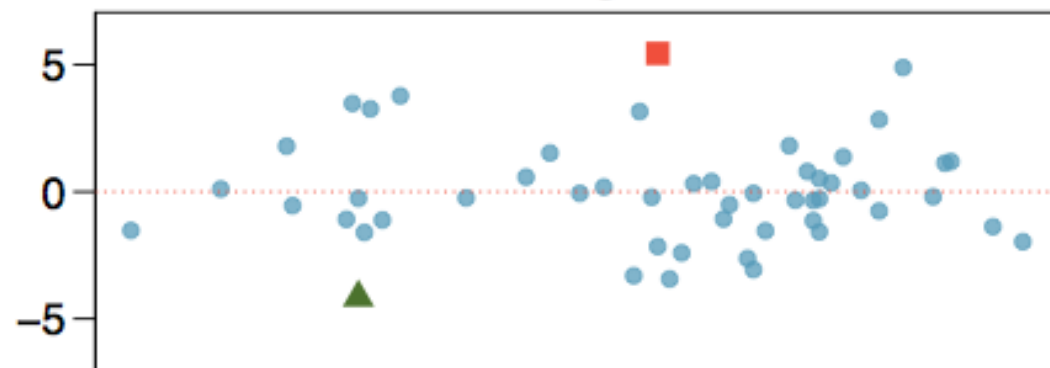


▲ RI:

$$\begin{aligned}\% \text{ HS grad} &= 81 & \% \text{ in poverty} &= 10.3 \\ \% \text{ in } \widehat{\text{poverty}} &= 64.68 - 0.62 * 81 = 14.46 \\ e &= \% \text{ in poverty} - \% \text{ in } \widehat{\text{poverty}} \\ &= 10.3 - 14.46 = -4.16\end{aligned}$$

■ DC:

$$\begin{aligned}\% \text{ HS grad} &= 86 & \% \text{ in poverty} &= 16.8 \\ \% \text{ in } \widehat{\text{poverty}} &= 64.68 - 0.62 * 86 = 11.36 \\ e &= \% \text{ in poverty} - \% \text{ in } \widehat{\text{poverty}} \\ &= 16.8 - 11.36 = 5.44\end{aligned}$$



Conditions: (2) Nearly normal residuals

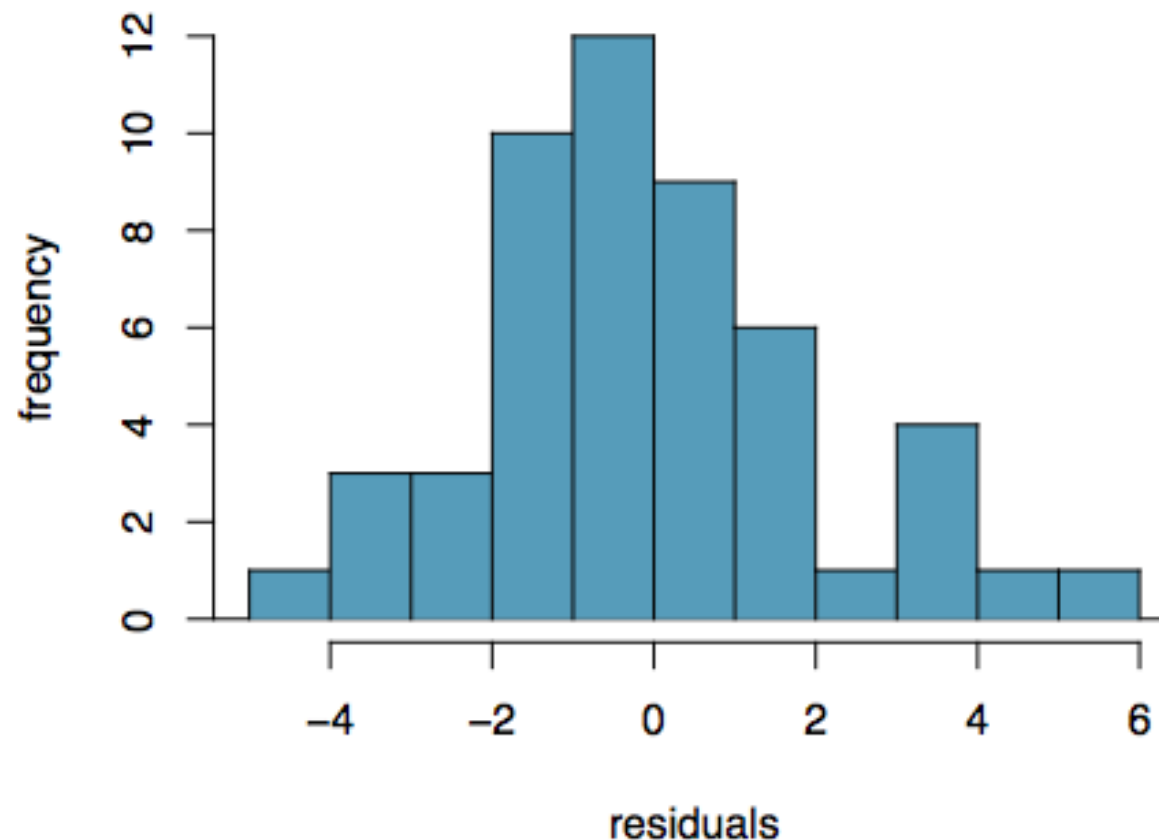
- The residuals should be nearly normal.

Conditions: (2) Nearly normal residuals

- The residuals should be nearly normal.
- This condition may not be satisfied when there are unusual observations that don't follow the trend of the rest of the data.

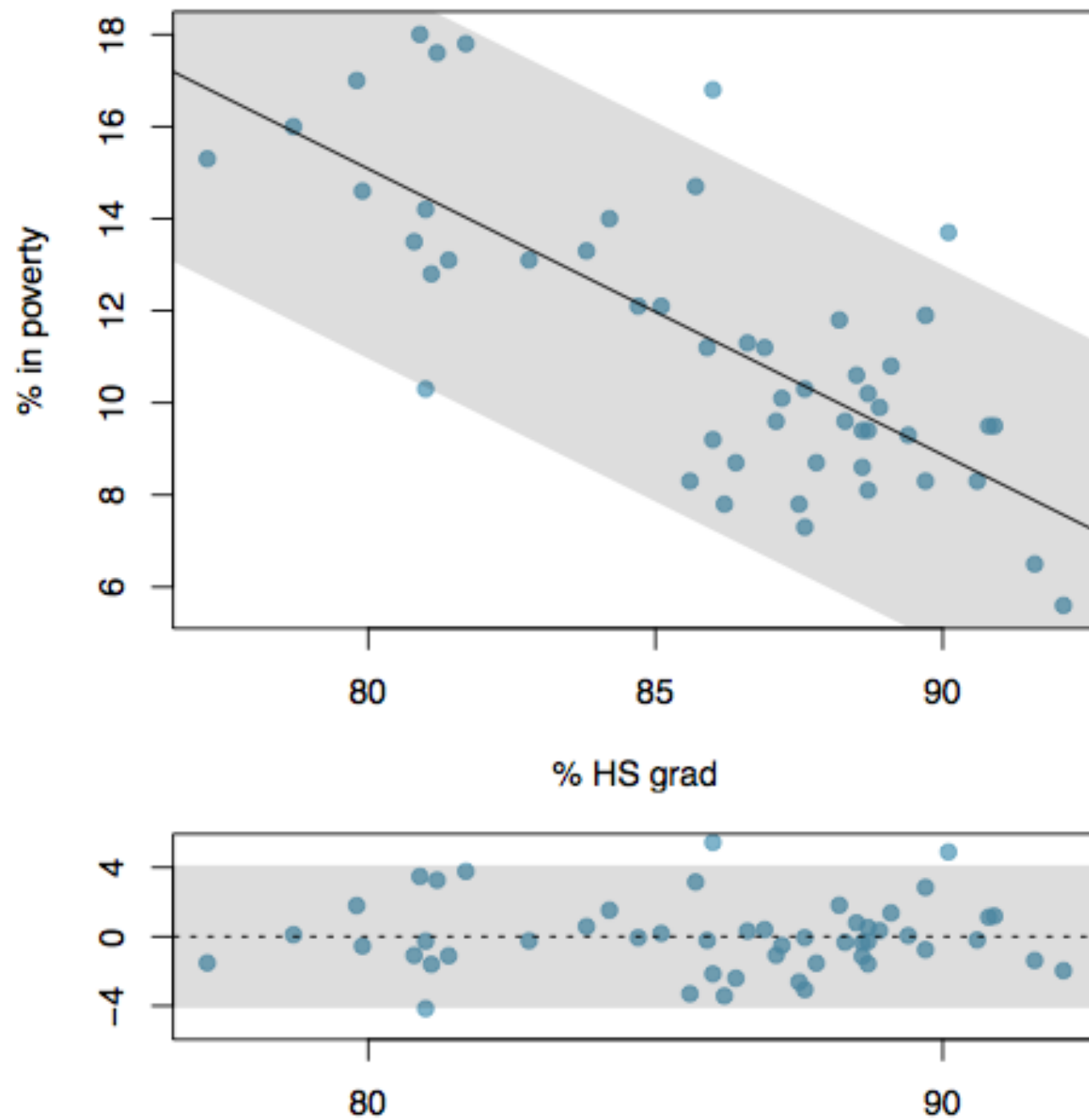
Conditions: (2) Nearly normal residuals

- The residuals should be nearly normal.
- This condition may not be satisfied when there are unusual observations that don't follow the trend of the rest of the data.
- Check using a histogram or normal probability plot of residuals.

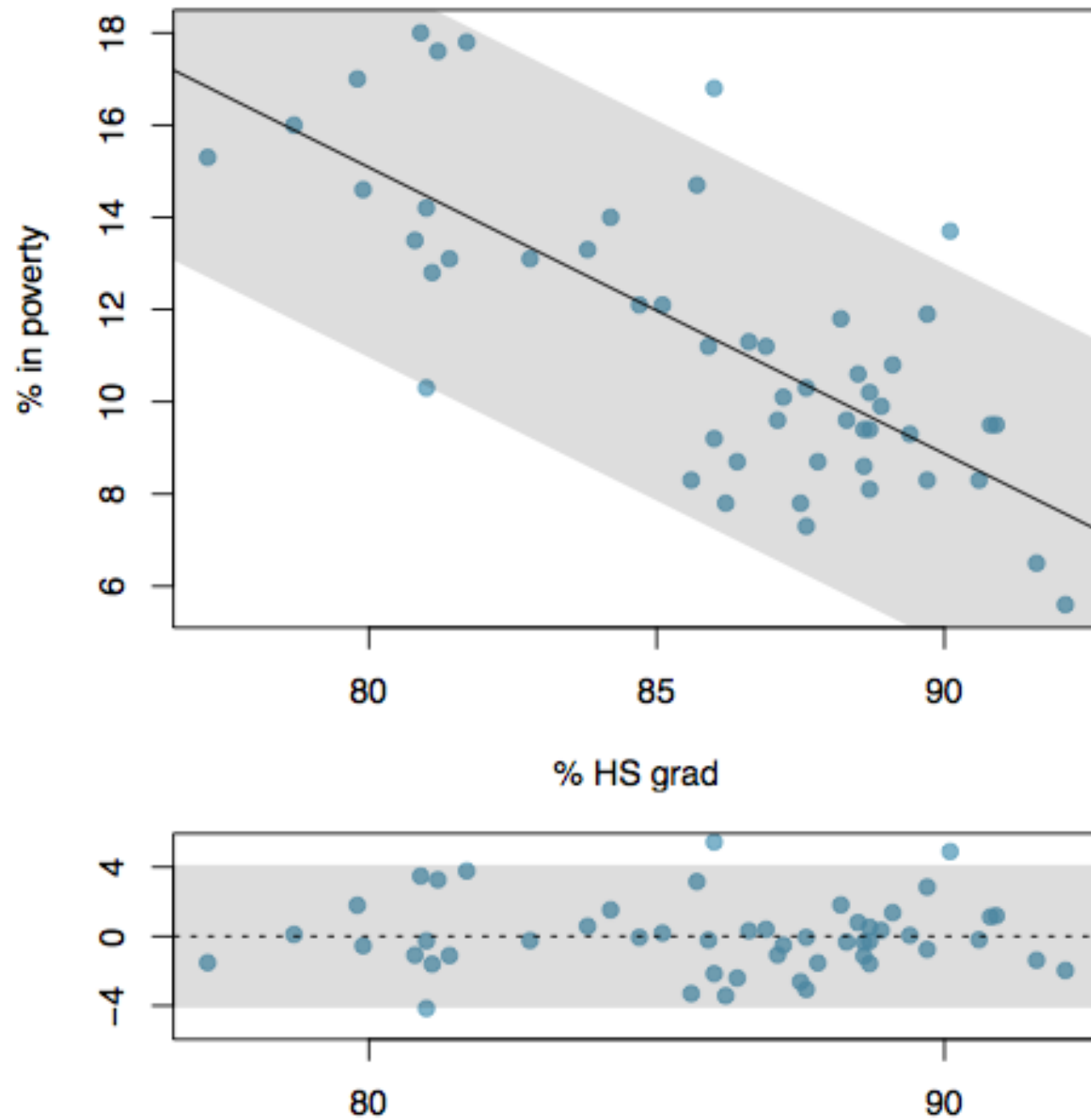


Conditions: (3) Constant variability

- The variability of points around the least squares line should be roughly constant.

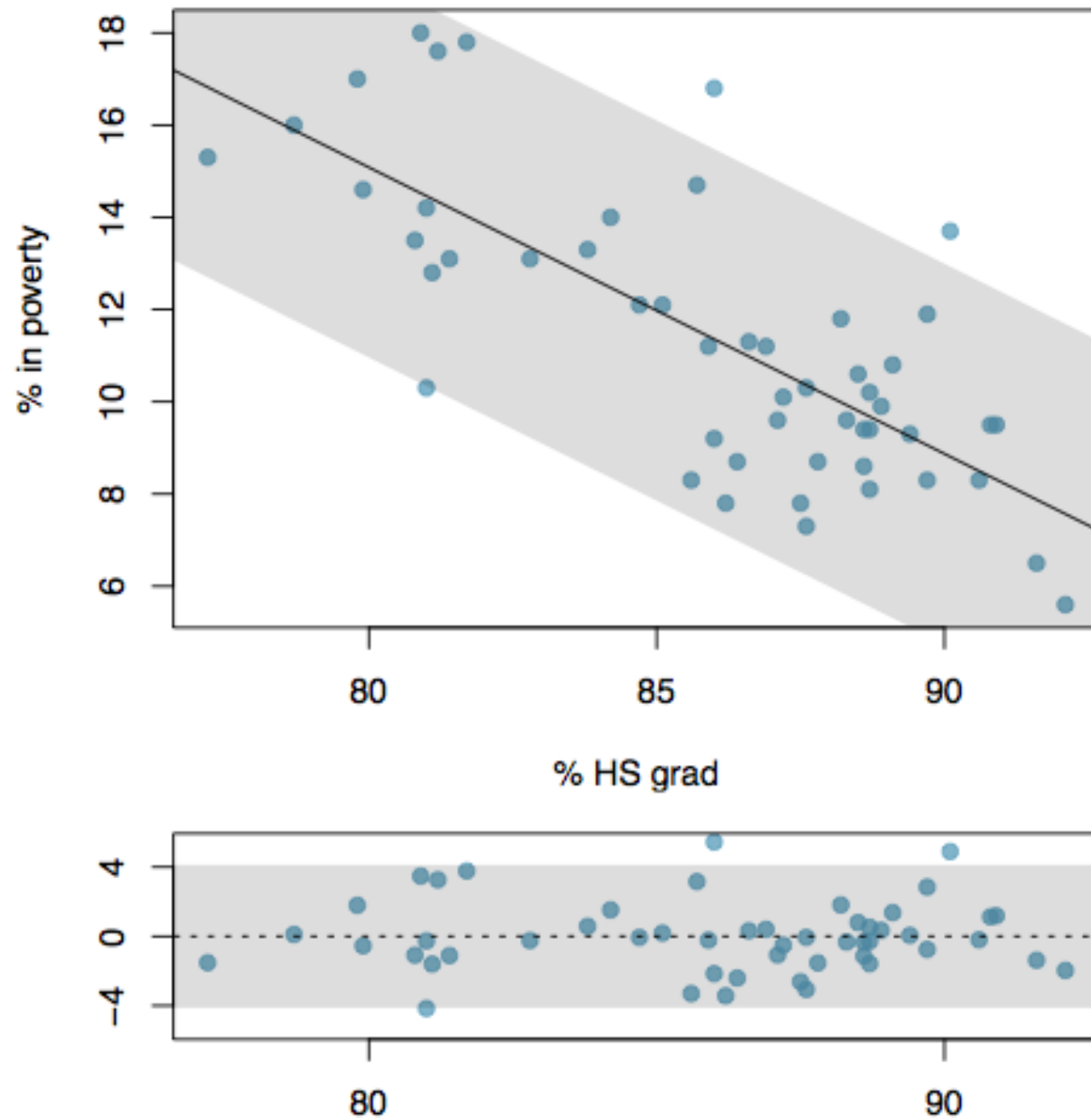


Conditions: (3) Constant variability



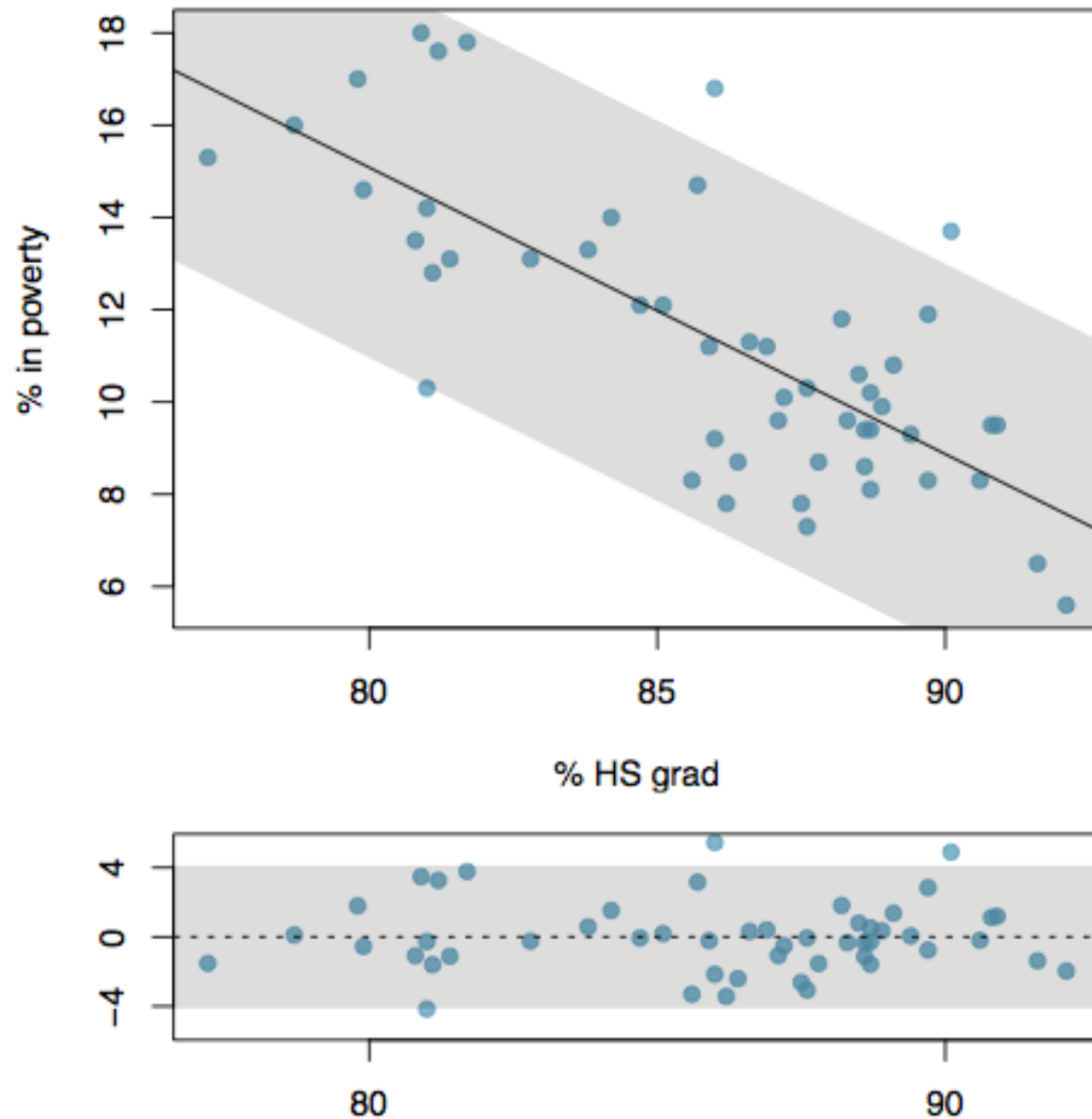
- The variability of points around the least squares line should be roughly constant.
- This implies that the variability of residuals around the 0 line should be roughly constant as well.

Conditions: (3) Constant variability



- The variability of points around the least squares line should be roughly constant.
- This implies that the variability of residuals around the 0 line should be roughly constant as well.
- Also called *homoscedasticity*.

Conditions: (3) Constant variability

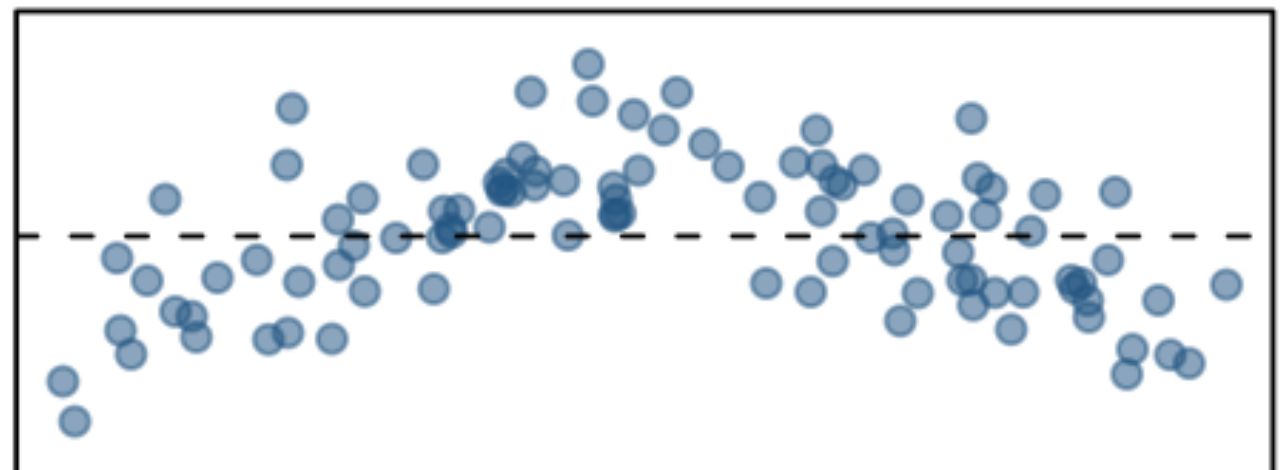
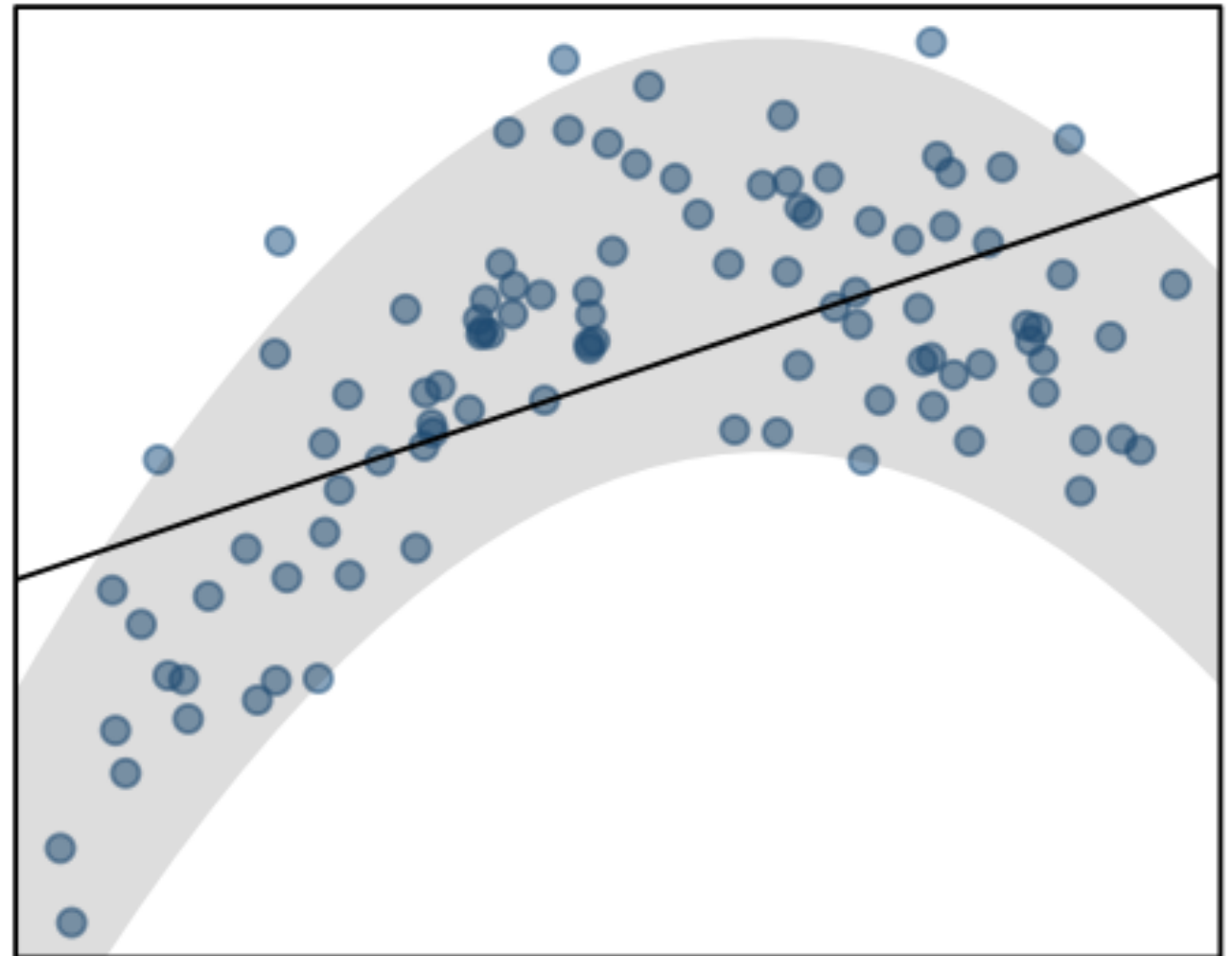


- The variability of points around the least squares line should be roughly constant.
- This implies that the variability of residuals around the 0 line should be roughly constant as well.
- Also called *homoscedasticity*.
- Check using a histogram or normal probability plot of residuals.

Checking conditions

What condition is this linear model obviously violating?

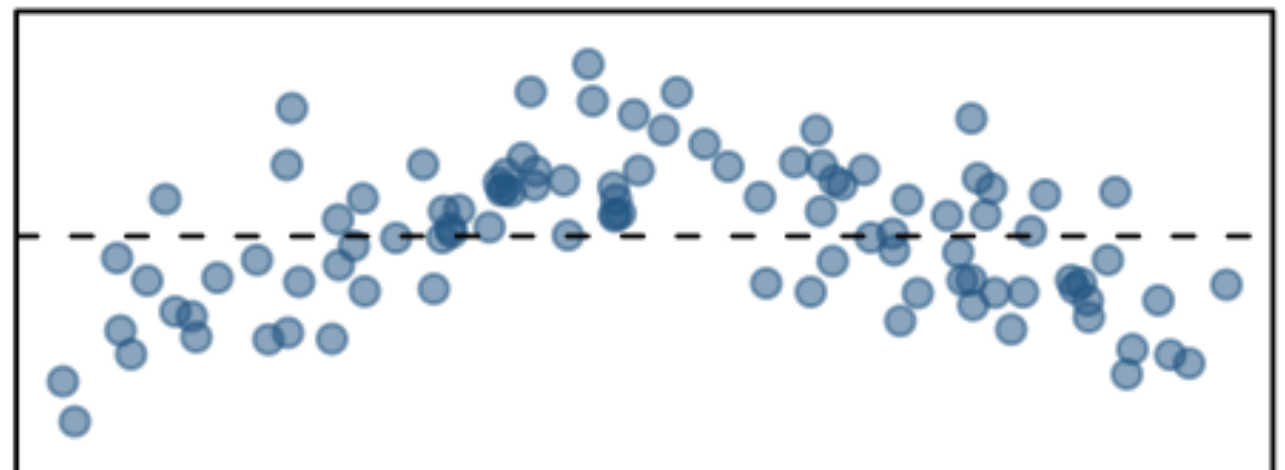
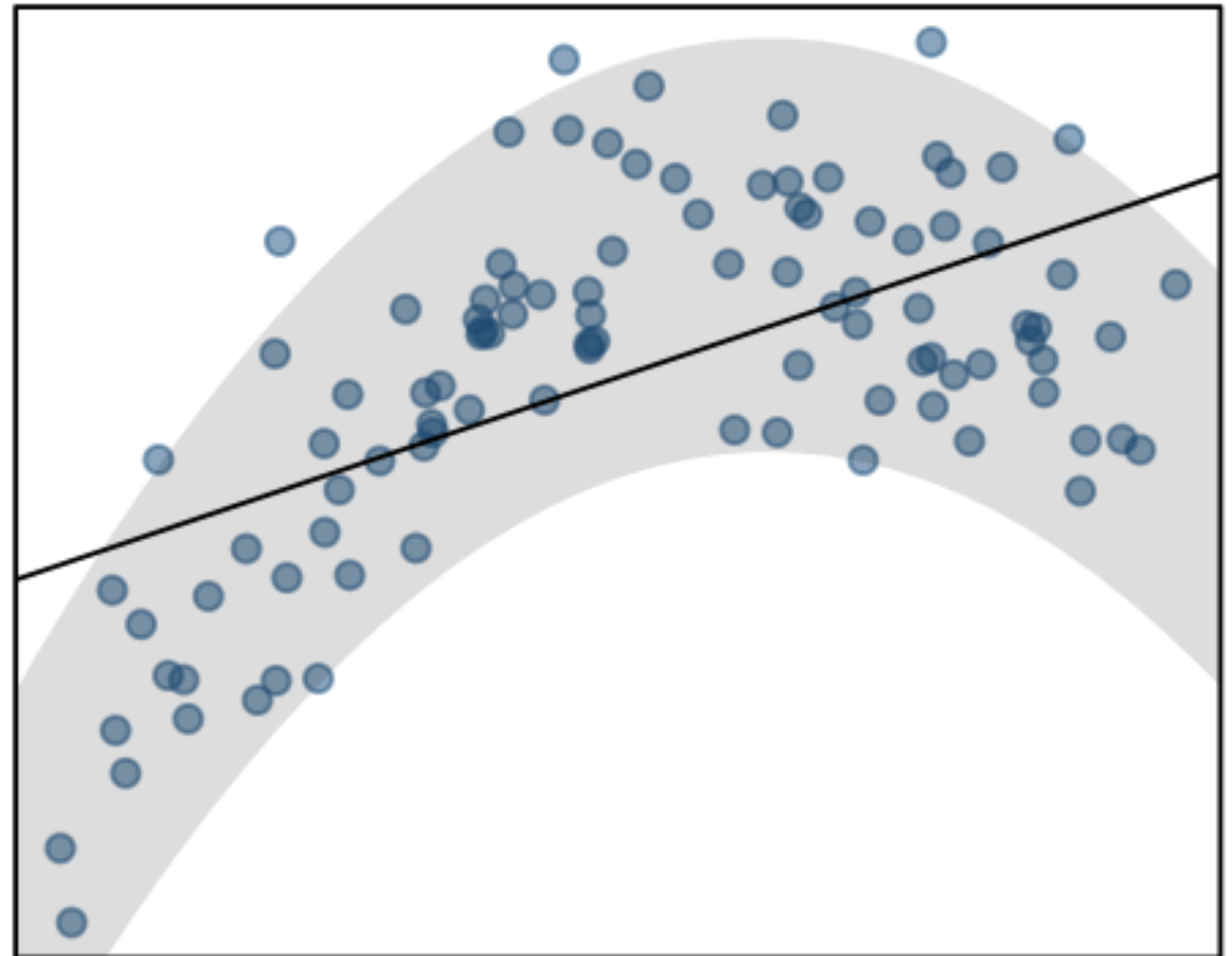
- (a) Constant variability
- (b) Linear relationship
- (c) Normal residuals
- (d) No extreme outliers



Checking conditions

What condition is this linear model obviously violating?

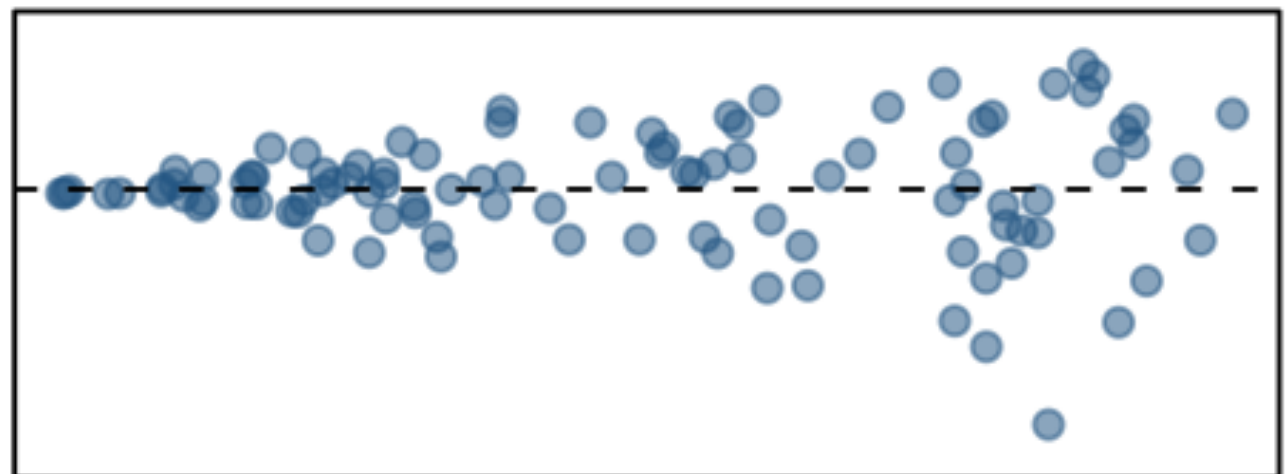
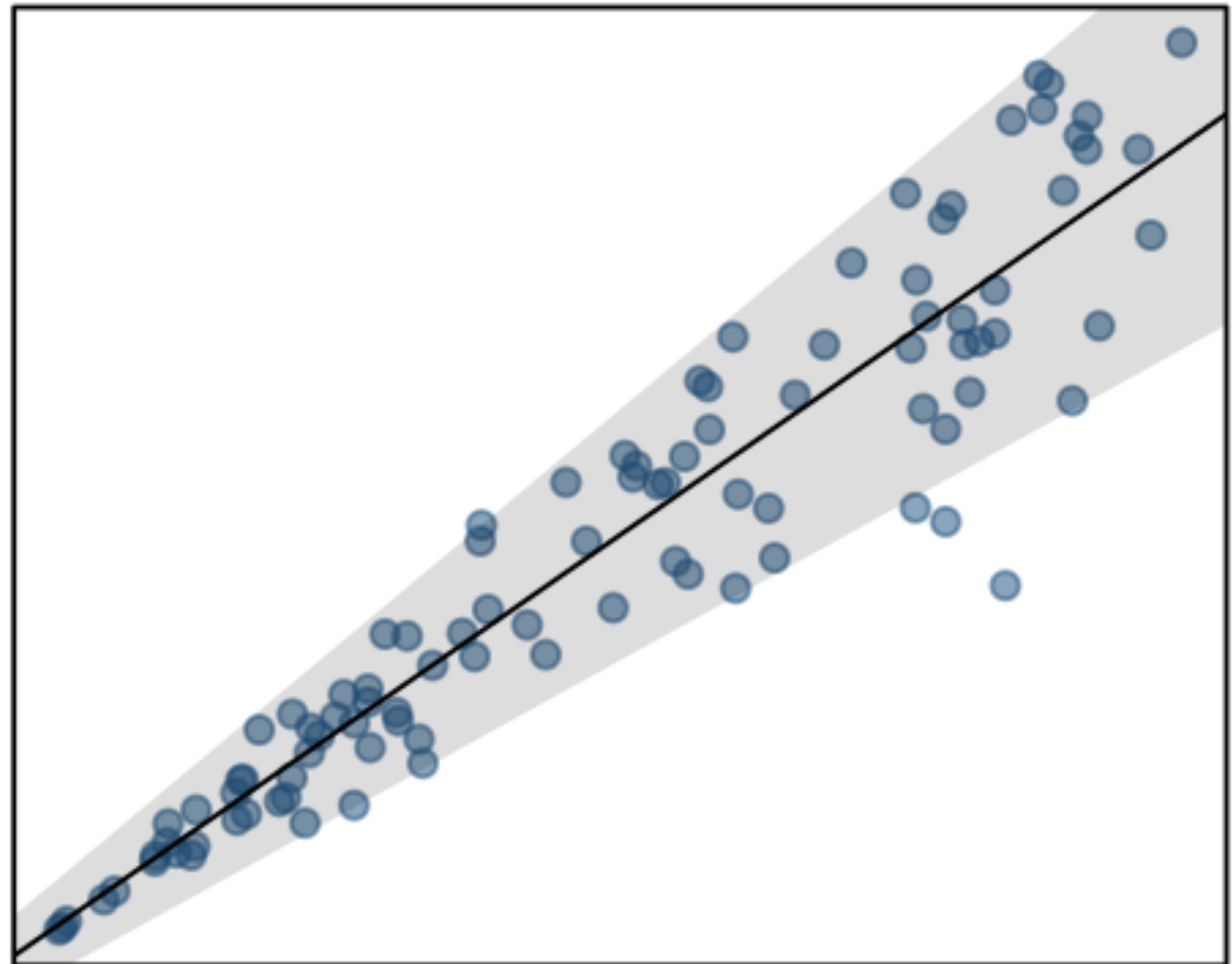
- (a) Constant variability
- (b) Linear relationship*
- (c) Normal residuals
- (d) No extreme outliers



Checking conditions

What condition is this linear model obviously violating?

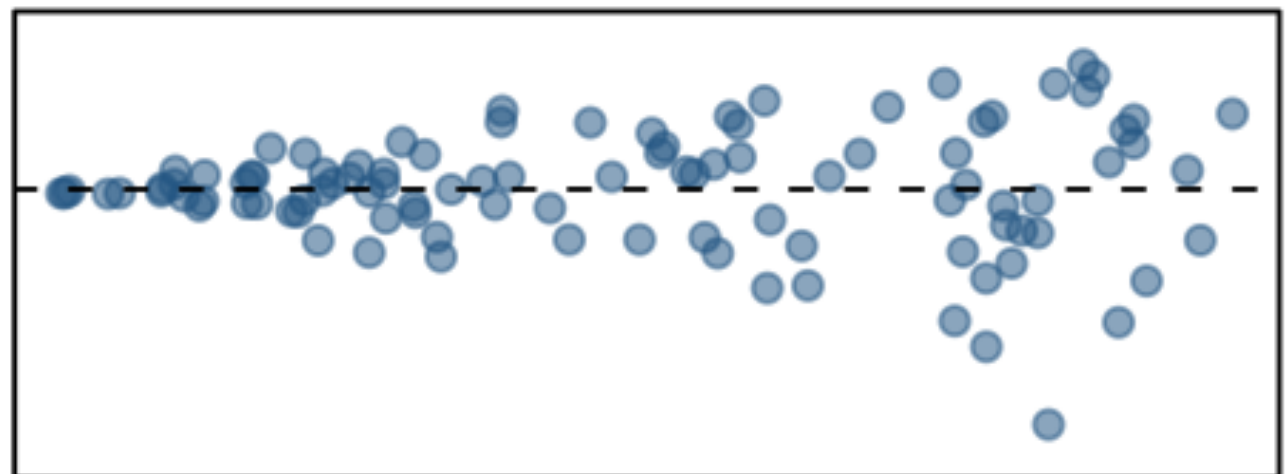
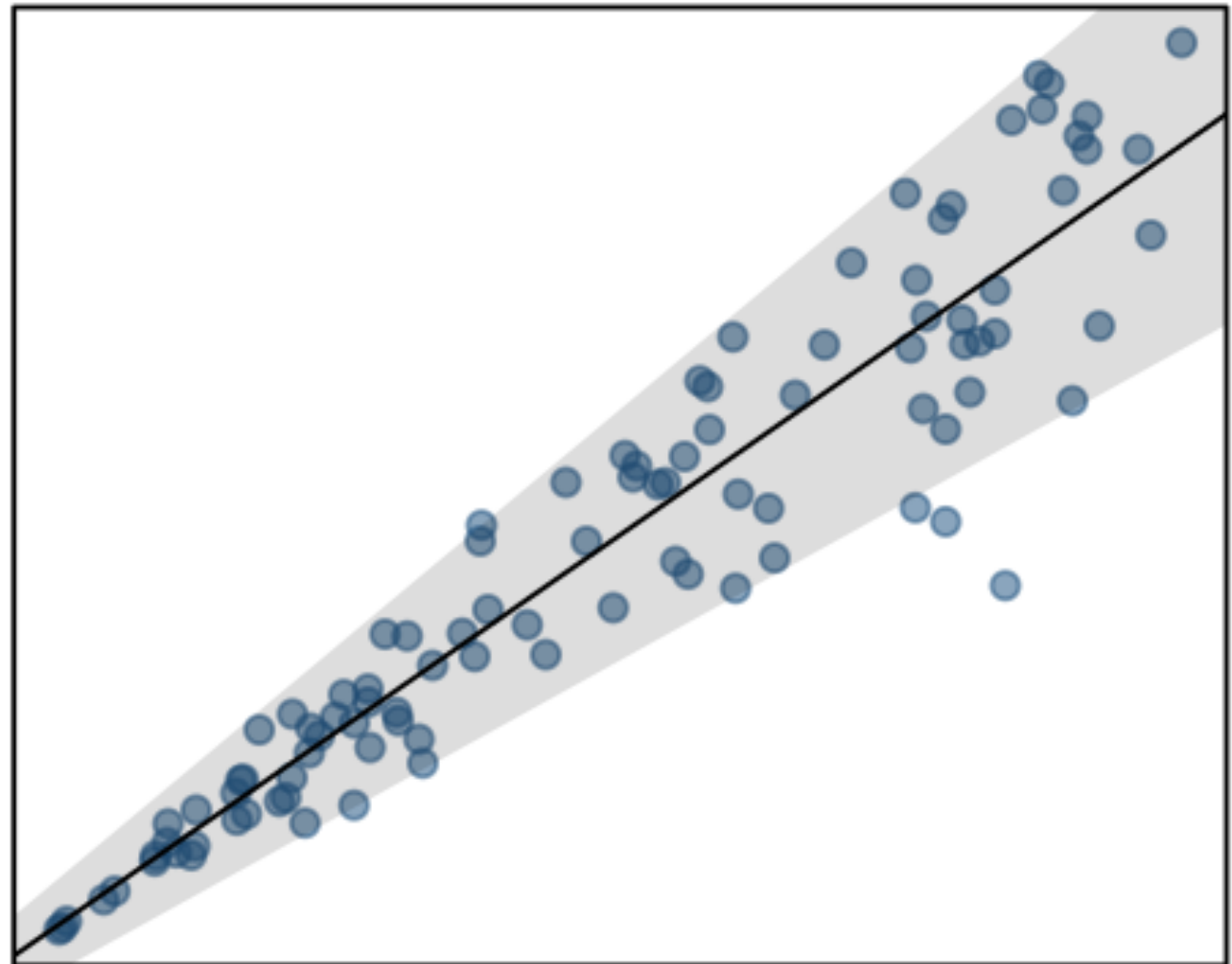
- (a) Constant variability
- (b) Linear relationship
- (c) Normal residuals
- (d) No extreme outliers



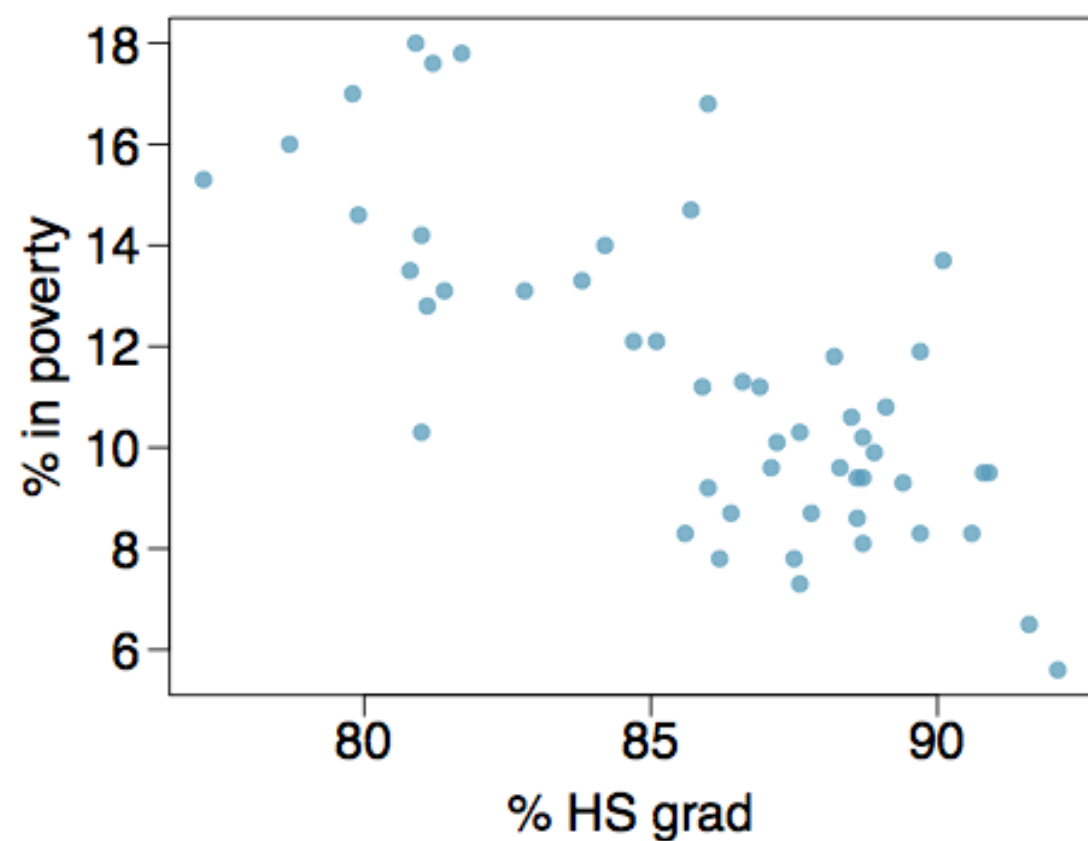
Checking conditions

What condition is this linear model obviously violating?

- (a) *Constant variability*
- (b) Linear relationship
- (c) Normal residuals
- (d) No extreme outliers



Given...



	% HS grad (x)	% in poverty (y)
mean	$\bar{x} = 86.01$	$\bar{y} = 11.35$
sd	$s_x = 3.73$	$s_y = 3.1$
correlation	$R = -0.75$	

Slope

The slope of the regression can be calculated as

$$b_1 = \frac{s_y}{s_x} R$$

Slope

The slope of the regression can be calculated as

$$b_1 = \frac{s_y}{s_x} R$$

In context...

$$b_1 = \frac{3.1}{3.73} \times -0.75 = -0.62$$

Slope

The slope of the regression can be calculated as

$$b_1 = \frac{s_y}{s_x} R$$

In context...

$$b_1 = \frac{3.1}{3.73} \times -0.75 = -0.62$$

Interpretation

For each additional % point in HS graduate rate, we would expect the % living in poverty to be lower on average by 0.62% points.

Intercept

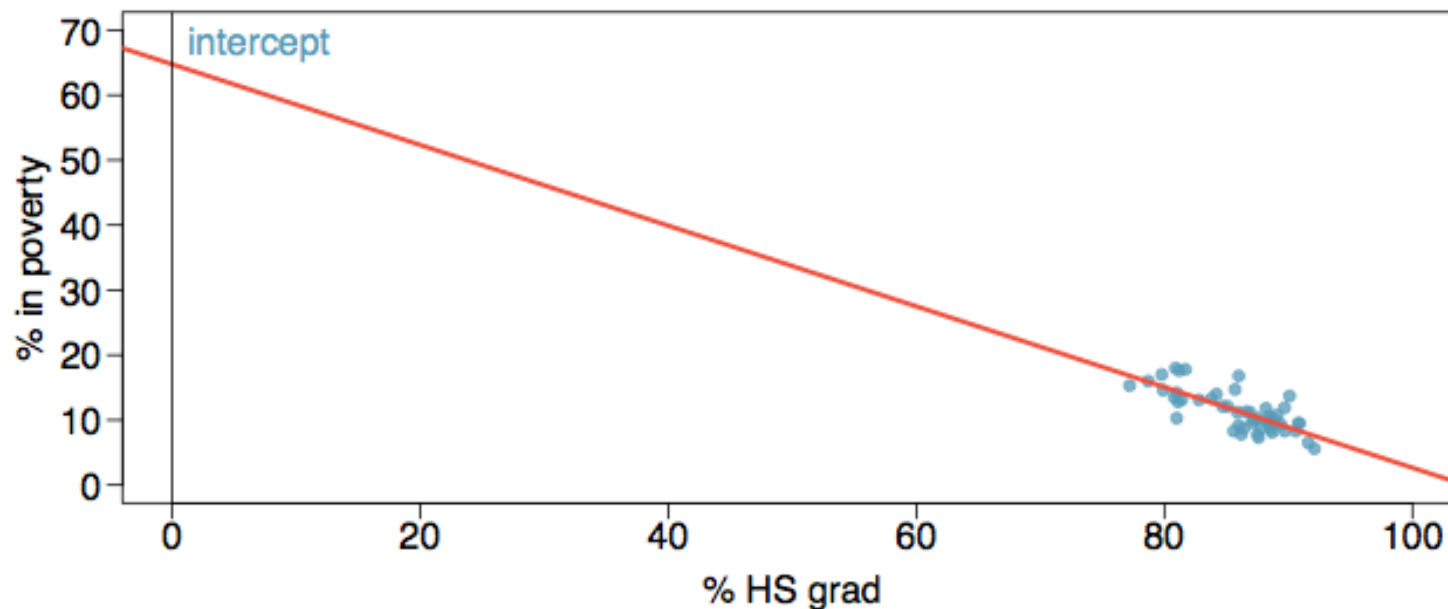
The intercept is where the regression line intersects the y-axis. The calculation of the intercept uses the fact that a regression line always passes through (\bar{x}, \bar{y}) .

$$b_0 = \bar{y} - b_1 \bar{x}$$

Intercept

The intercept is where the regression line intersects the y-axis. The calculation of the intercept uses the fact that a regression line always passes through (\bar{x}, \bar{y}) .

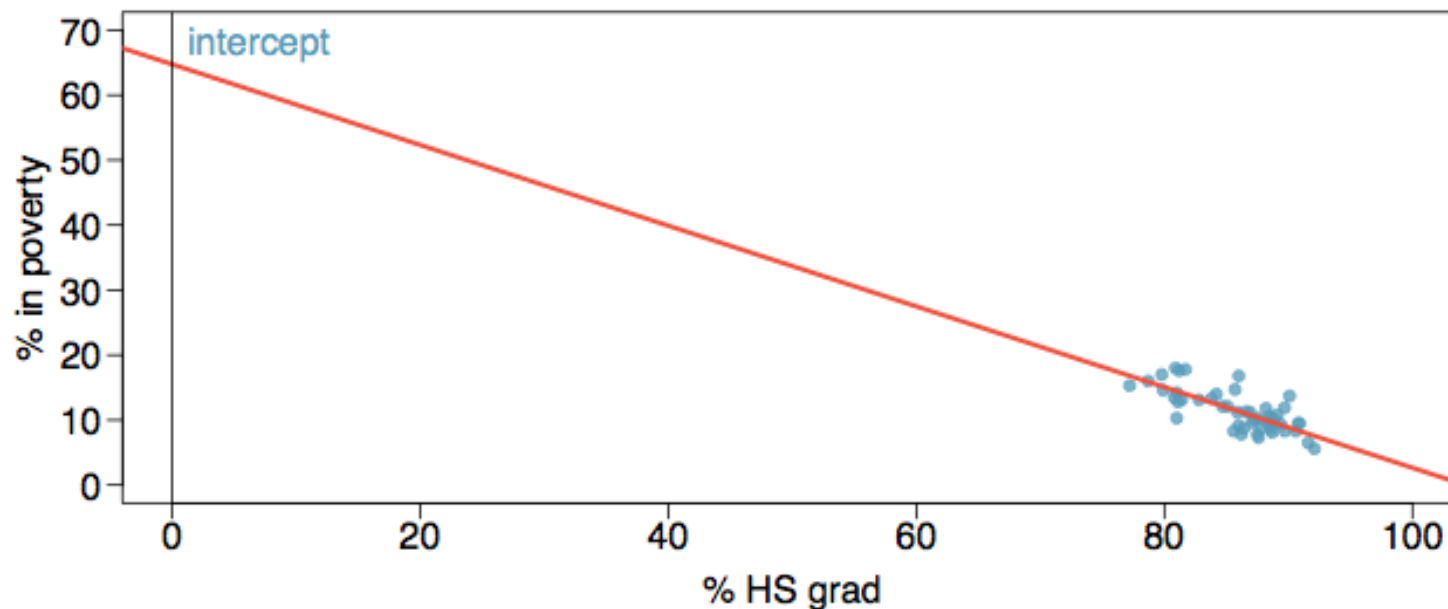
$$b_0 = \bar{y} - b_1 \bar{x}$$



Intercept

The intercept is where the regression line intersects the y-axis. The calculation of the intercept uses the fact that a regression line always passes through (\bar{x}, \bar{y}) .

$$b_0 = \bar{y} - b_1 \bar{x}$$



$$\begin{aligned} b_0 &= 11.35 - (-0.62) \times 86.01 \\ &= 64.68 \end{aligned}$$

Which of the following is the correct interpretation of the intercept?

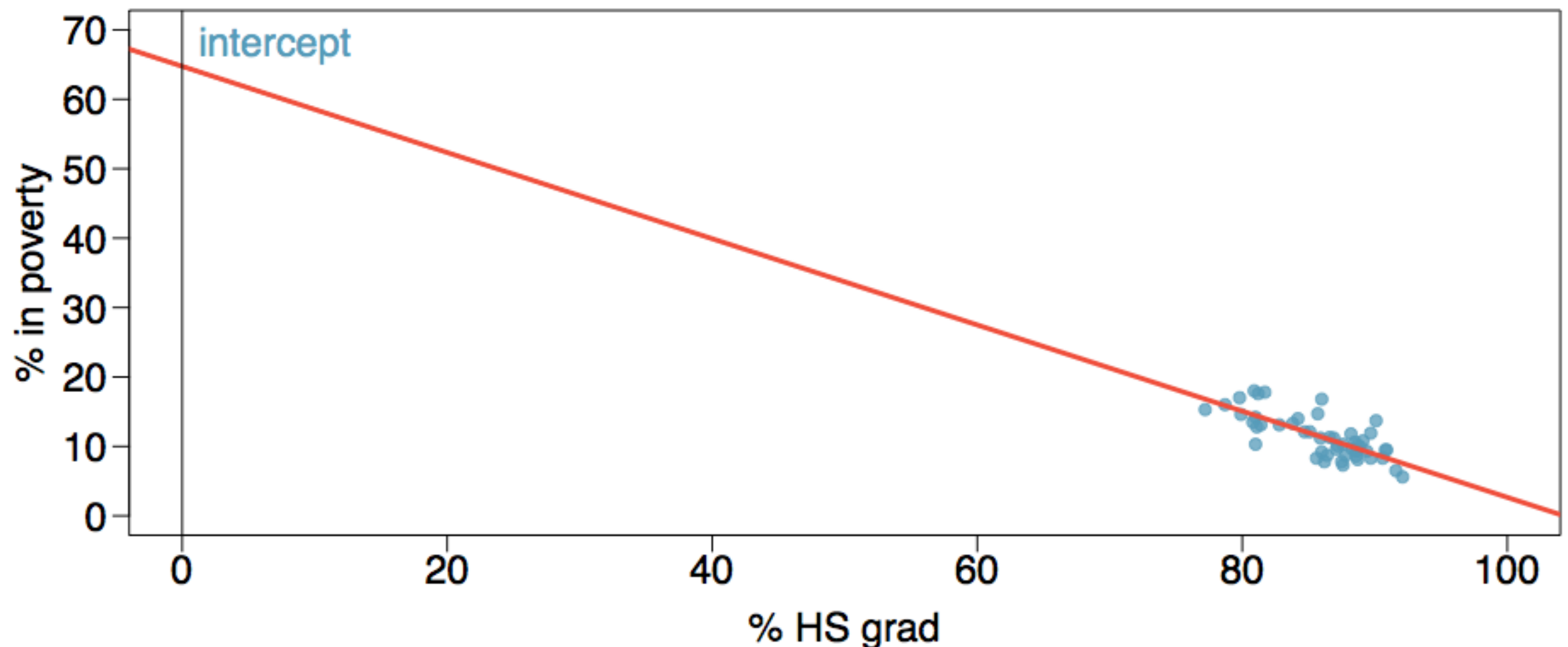
- (a) For each % point increase in HS graduate rate, % living in poverty is expected to increase on average by 64.68%.
- (b) For each % point decrease in HS graduate rate, % living in poverty is expected to increase on average by 64.68%.
- (c) Having no HS graduates leads to 64.68% of residents living below the poverty line.
- (d) States with no HS graduates are expected on average to have 64.68% of residents living below the poverty line.
- (e) In states with no HS graduates % living in poverty is expected to increase on average by 64.68%.

Which of the following is the correct interpretation of the intercept?

- (a) For each % point increase in HS graduate rate, % living in poverty is expected to increase on average by 64.68%.
- (b) For each % point decrease in HS graduate rate, % living in poverty is expected to increase on average by 64.68%.
- (c) Having no HS graduates leads to 64.68% of residents living below the poverty line.
- (d) States with no HS graduates are expected on average to have 64.68% of residents living below the poverty line.*
- (e) In states with no HS graduates % living in poverty is expected to increase on average by 64.68%.

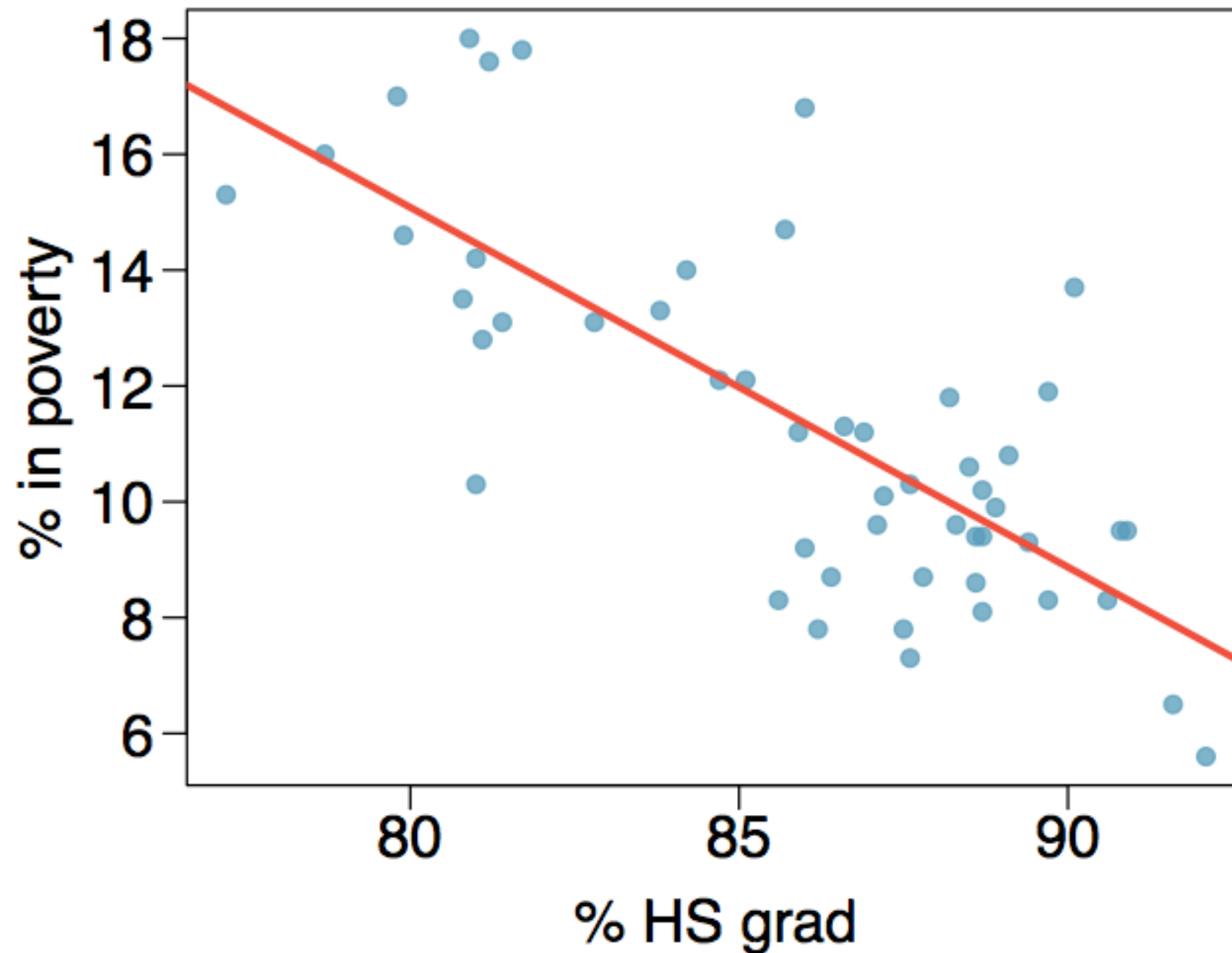
More on the intercept

Since there are no states in the dataset with no HS graduates, the intercept is of no interest, not very useful, and also not reliable since the predicted value of the intercept is so far from the bulk of the data.



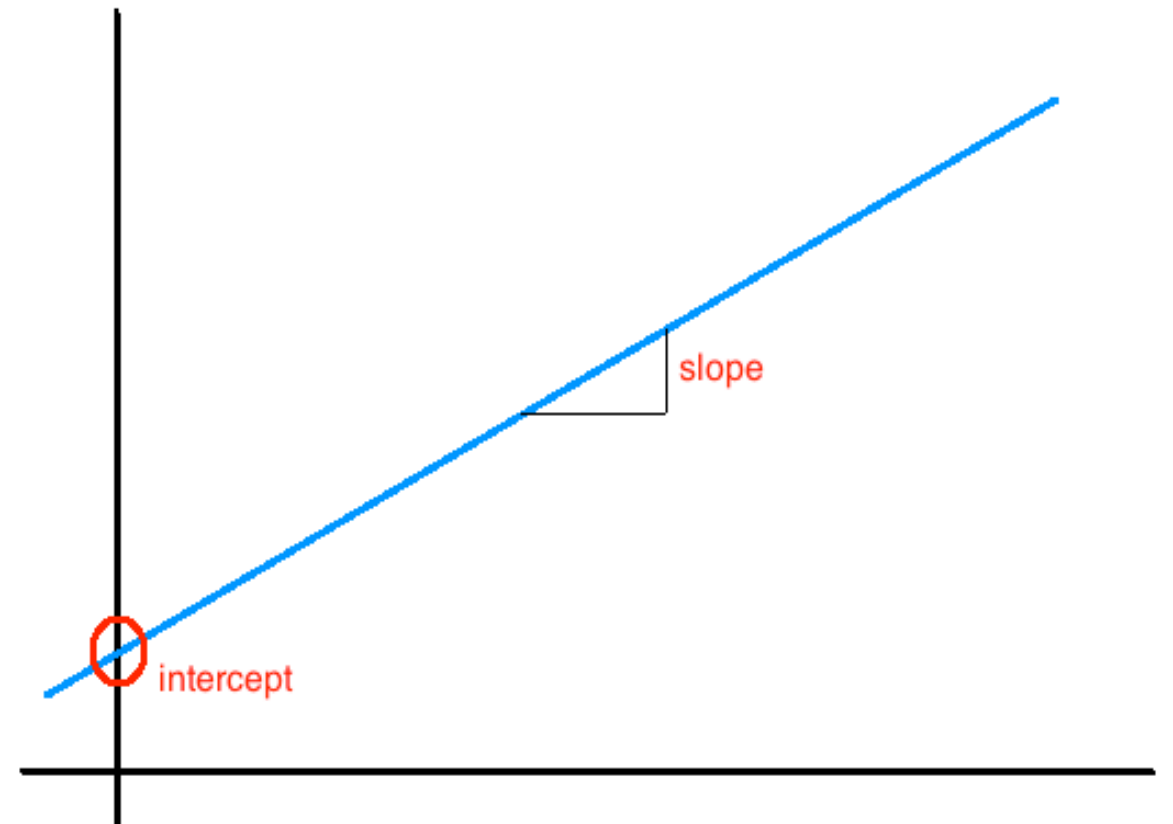
Regression line

$$\widehat{\% \text{ in poverty}} = 64.68 - 0.62 \% \text{ HS grad}$$



Interpretation of slope and intercept

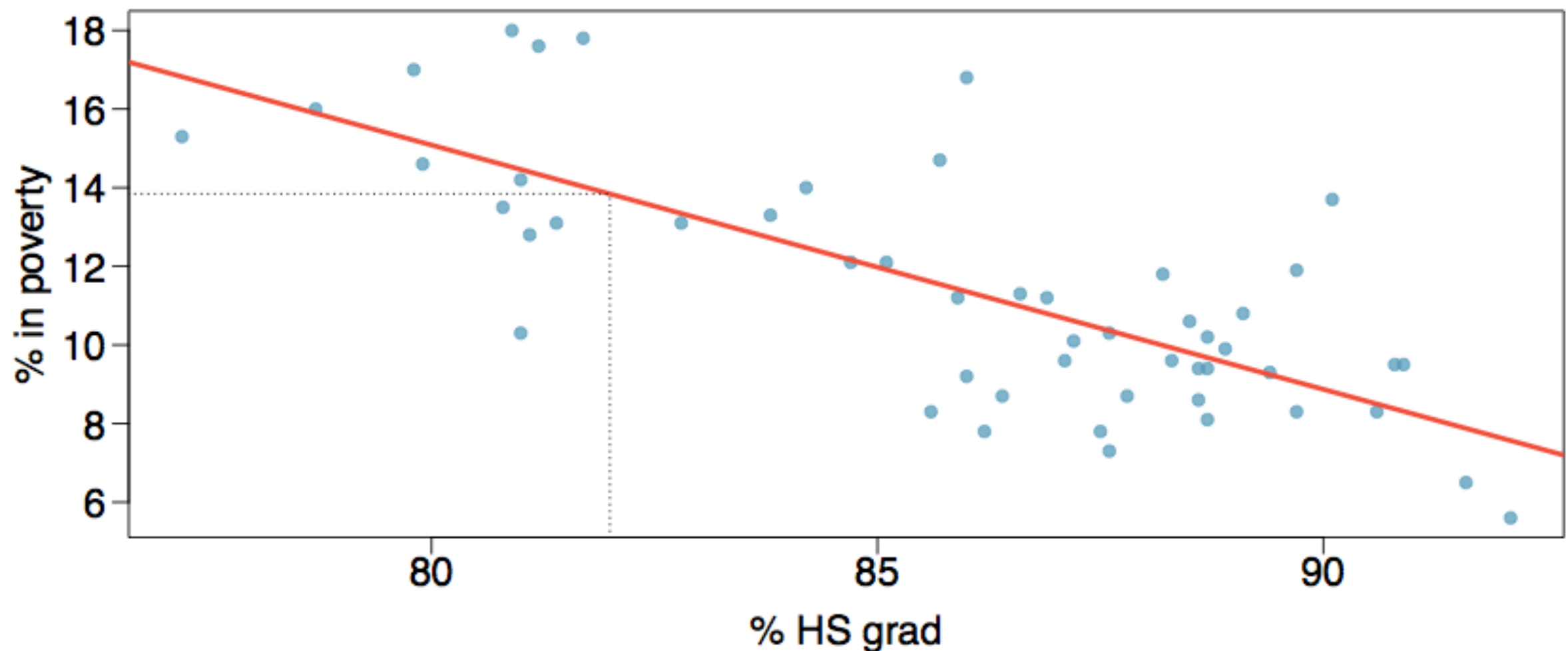
- *Intercept*: When $x = 0$, y is expected to equal the intercept.
- *Slope*: For each unit in x , y is expected to increase / decrease on average by the slope.



Note: These statements are not causal, unless the study is a randomized controlled experiment.

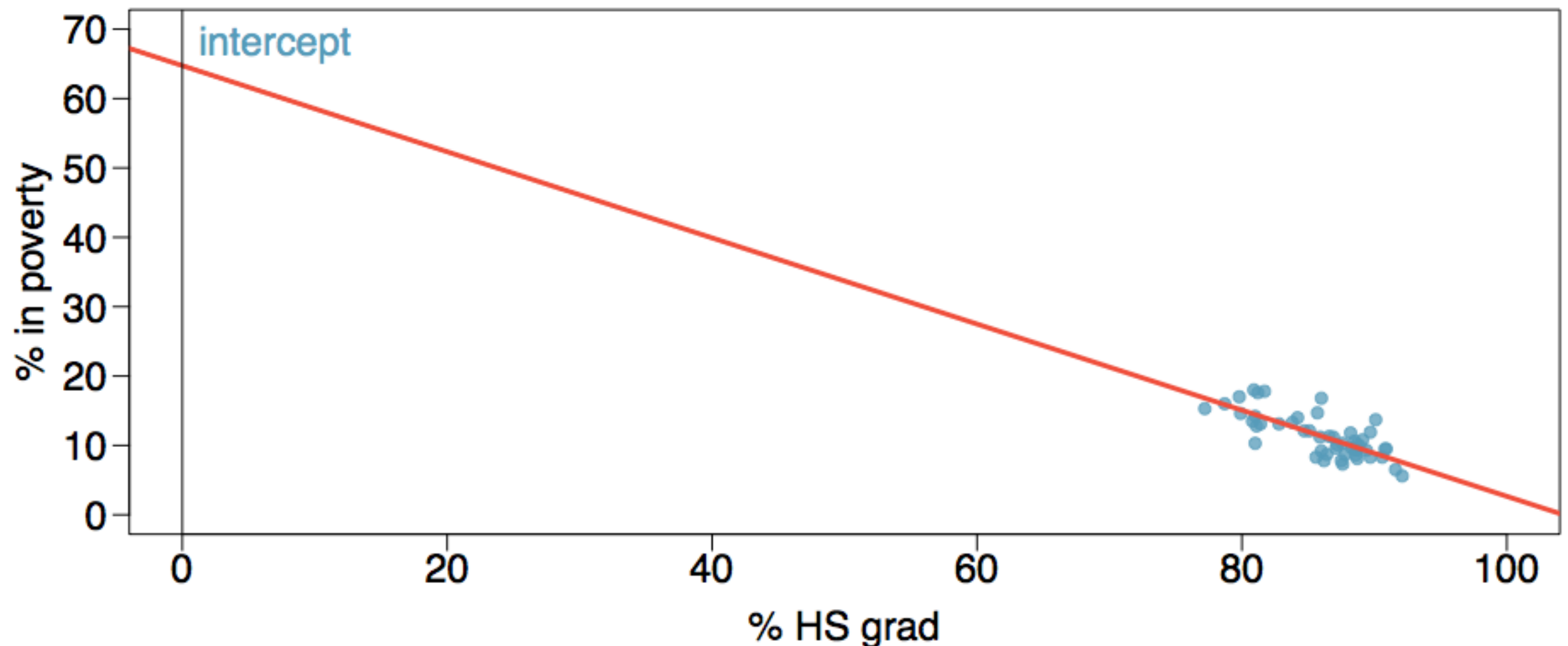
Prediction

- Using the linear model to predict the value of the response variable for a given value of the explanatory variable is called *prediction*, simply by plugging in the value of x in the linear model equation.
- There will be some uncertainty associated with the predicted value.



Extrapolation

- Applying a model estimate to values outside of the realm of the original data is called *extrapolation*.
- Sometimes the intercept might be an extrapolation.



Examples of extrapolation

Momentous sprint at the 2156 Olympics?

Women sprinters are closing the gap on men and may one day overtake them.

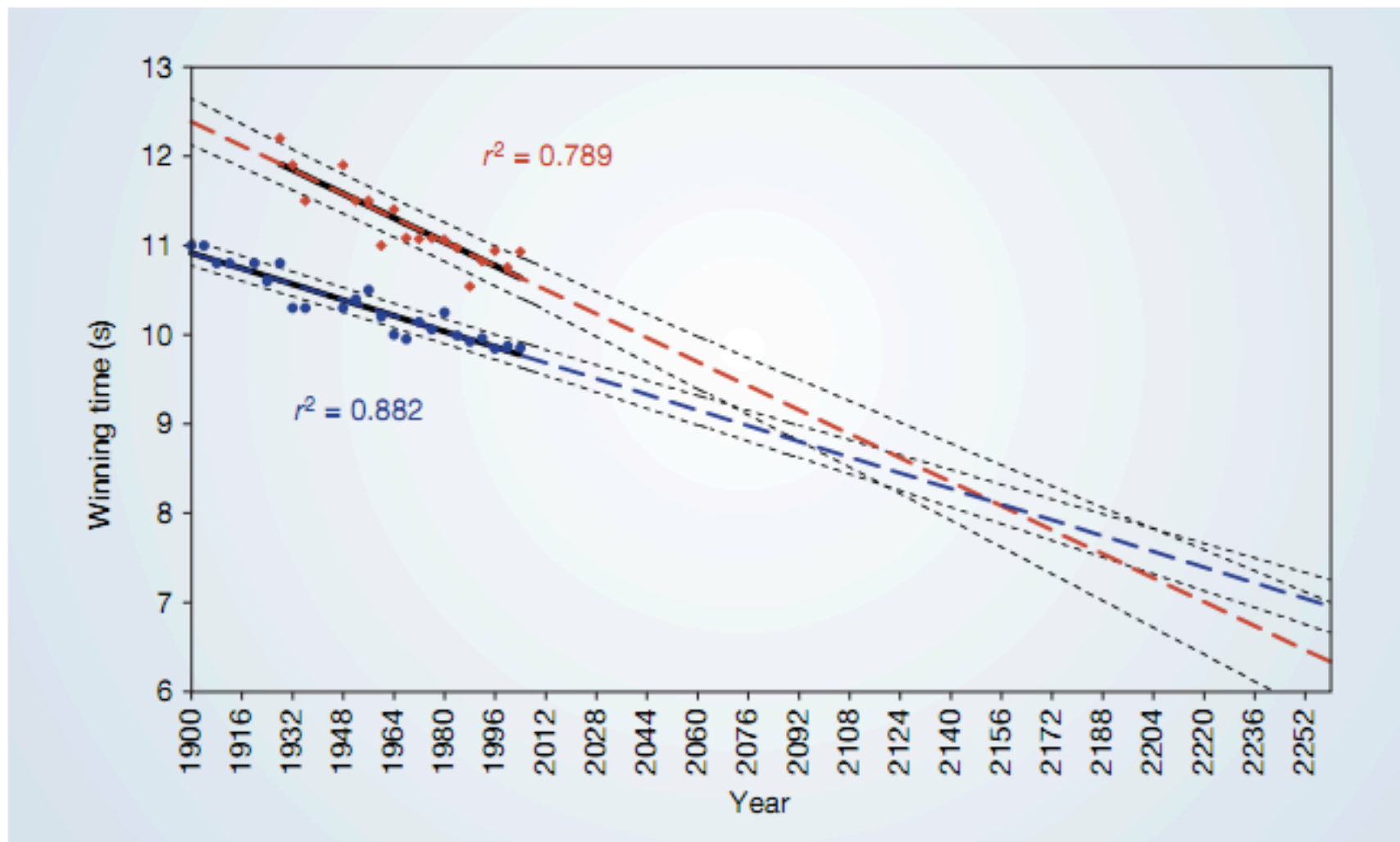


Figure 1 The winning Olympic 100-metre sprint times for men (blue points) and women (red points), with superimposed best-fit linear regression lines (solid black lines) and coefficients of determination. The regression lines are extrapolated (broken blue and red lines for men and women, respectively) and 95% confidence intervals (dotted black lines) based on the available points are superimposed. The projections intersect just before the 2156 Olympics, when the winning women's 100-metre sprint time of 8.079 s will be faster than the men's at 8.098 s.

R^2

- The strength of the fit of a linear model is most commonly evaluated using R^2 .

R^2

- The strength of the fit of a linear model is most commonly evaluated using R^2 .
- R^2 is calculated as the square of the correlation coefficient.

R^2

- The strength of the fit of a linear model is most commonly evaluated using R^2 .
- R^2 is calculated as the square of the correlation coefficient.
- It tells us what percent of variability in the response variable is explained by the model.

R^2

- The strength of the fit of a linear model is most commonly evaluated using R^2 .
- R^2 is calculated as the square of the correlation coefficient.
- It tells us what percent of variability in the response variable is explained by the model.
- The remainder of the variability is explained by variables not included in the model or by inherent randomness in the data.

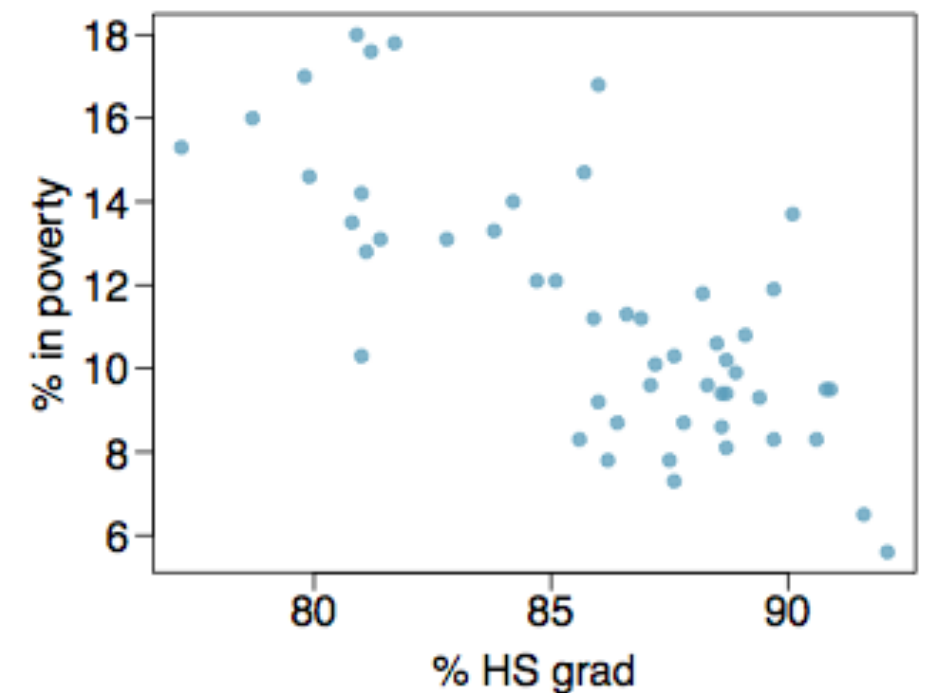
R^2

- The strength of the fit of a linear model is most commonly evaluated using R^2 .
- R^2 is calculated as the square of the correlation coefficient.
- It tells us what percent of variability in the response variable is explained by the model.
- The remainder of the variability is explained by variables not included in the model or by inherent randomness in the data.
- For the model we've been working with, $R^2 = -0.62^2 = 0.38$.

Interpretation of R^2

Which of the below is the correct interpretation of $R = -0.62$, $R^2 = 0.38$?

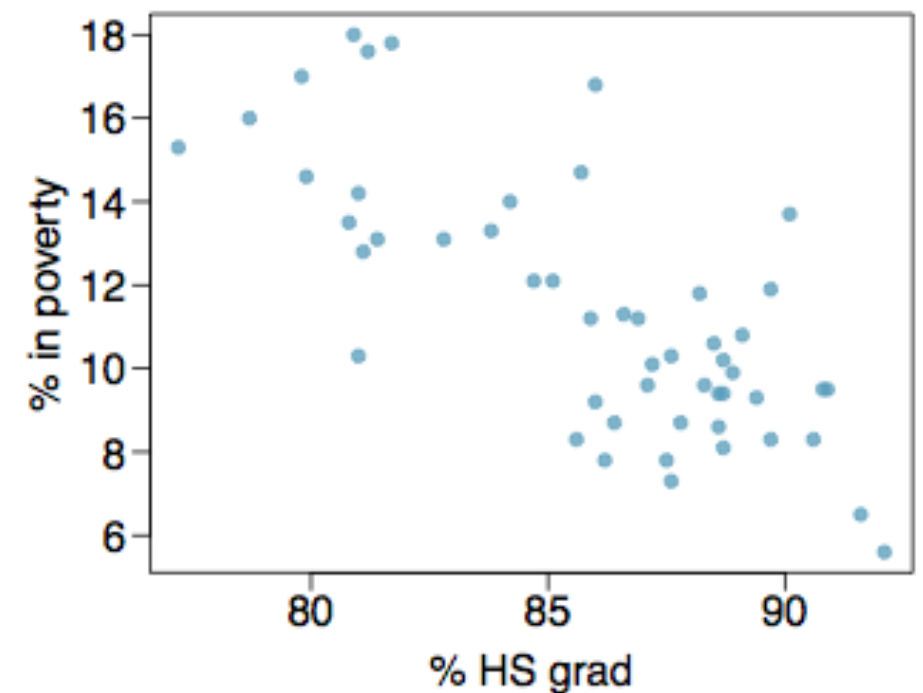
- (a) 38% of the variability in the % of HG graduates among the 51 states is explained by the model.
- (b) 38% of the variability in the % of residents living in poverty among the 51 states is explained by the model.
- (c) 38% of the time % HS graduates predict % living in poverty correctly.
- (d) 62% of the variability in the % of residents living in poverty among the 51 states is explained by the model.



Interpretation of R^2

Which of the below is the correct interpretation of $R = -0.62$, $R^2 = 0.38$?

- (a) 38% of the variability in the % of HG graduates among the 51 states is explained by the model.
- (b) 38% of the variability in the % of residents living in poverty among the 51 states is explained by the model.*
- (c) 38% of the time % HS graduates predict % living in poverty correctly.
- (d) 62% of the variability in the % of residents living in poverty among the 51 states is explained by the model.



Poverty vs. region (east, west)

$$\widehat{poverty} = 11.17 + 0.38 \times west$$

- Explanatory variable: region, *reference level*: east
- *Intercept*: The estimated average poverty percentage in eastern states is 11.17%

Poverty vs. region (east, west)

$$\widehat{poverty} = 11.17 + 0.38 \times west$$

- Explanatory variable: region, *reference level*: east
- *Intercept*: The estimated average poverty percentage in eastern states is 11.17%
 - This is the value we get if we plug in 0 for the explanatory variable

Poverty vs. region (east, west)

$$\widehat{poverty} = 11.17 + 0.38 \times west$$

- Explanatory variable: region, *reference level*: east
- *Intercept*: The estimated average poverty percentage in eastern states is 11.17%
 - This is the value we get if we plug in 0 for the explanatory variable
- *Slope*: The estimated average poverty percentage in western states is 0.38% higher than eastern states.

Poverty vs. region (east, west)

$$\widehat{poverty} = 11.17 + 0.38 \times west$$

- Explanatory variable: region, *reference level*: east
- *Intercept*: The estimated average poverty percentage in eastern states is 11.17%
 - This is the value we get if we plug in 0 for the explanatory variable
- *Slope*: The estimated average poverty percentage in western states is 0.38% higher than eastern states.
 - Then, the estimated average poverty percentage in western states is $11.17 + 0.38 = 11.55\%$.

Poverty vs. region (east, west)

$$\widehat{poverty} = 11.17 + 0.38 \times west$$

- Explanatory variable: region, *reference level*: east
- *Intercept*: The estimated average poverty percentage in eastern states is 11.17%
 - This is the value we get if we plug in 0 for the explanatory variable
- *Slope*: The estimated average poverty percentage in western states is 0.38% higher than eastern states.
 - Then, the estimated average poverty percentage in western states is $11.17 + 0.38 = 11.55\%$.
 - This is the value we get if we plug in 1 for the explanatory variable

Poverty vs. region (northeast, midwest, west, south)

Which region (northeast, midwest, west, or south) is the reference level?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.50	0.87	10.94	0.00
region4midwest	0.03	1.15	0.02	0.98
region4west	1.79	1.13	1.59	0.12
region4south	4.16	1.07	3.87	0.00

- (a) northeast
- (b) midwest
- (c) west
- (d) south
- (e) cannot tell

Poverty vs. region (northeast, midwest, west, south)

Which region (northeast, midwest, west, or south) is the reference level?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.50	0.87	10.94	0.00
region4midwest	0.03	1.15	0.02	0.98
region4west	1.79	1.13	1.59	0.12
region4south	4.16	1.07	3.87	0.00

(a) northeast

(b) midwest

(c) west

(d) south

(e) cannot tell

Poverty vs. region (northeast, midwest, west, south)

Which region (northeast, midwest, west, or south) has the lowest poverty percentage?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.50	0.87	10.94	0.00
region4midwest	0.03	1.15	0.02	0.98
region4west	1.79	1.13	1.59	0.12
region4south	4.16	1.07	3.87	0.00

- (a) northeast
- (b) midwest
- (c) west
- (d) south
- (e) cannot tell

Poverty vs. region (northeast, midwest, west, south)

Which region (northeast, midwest, west, or south) has the lowest poverty percentage?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.50	0.87	10.94	0.00
region4midwest	0.03	1.15	0.02	0.98
region4west	1.79	1.13	1.59	0.12
region4south	4.16	1.07	3.87	0.00

(a) northeast

(b) midwest

(c) west

(d) south

(e) cannot tell

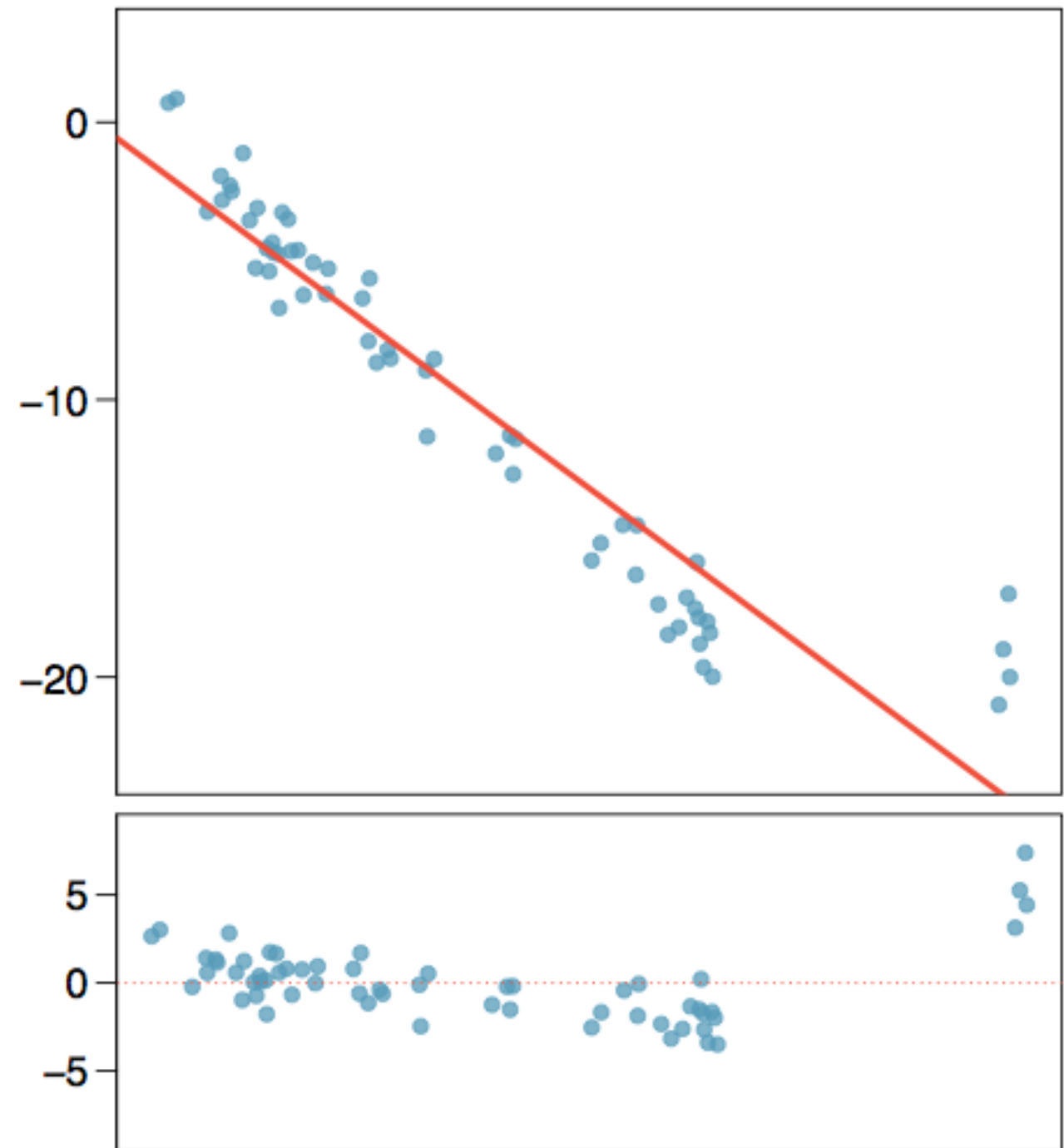
Types of outliers in linear regression

Types of outliers

How do outliers influence the least squares line in this plot?

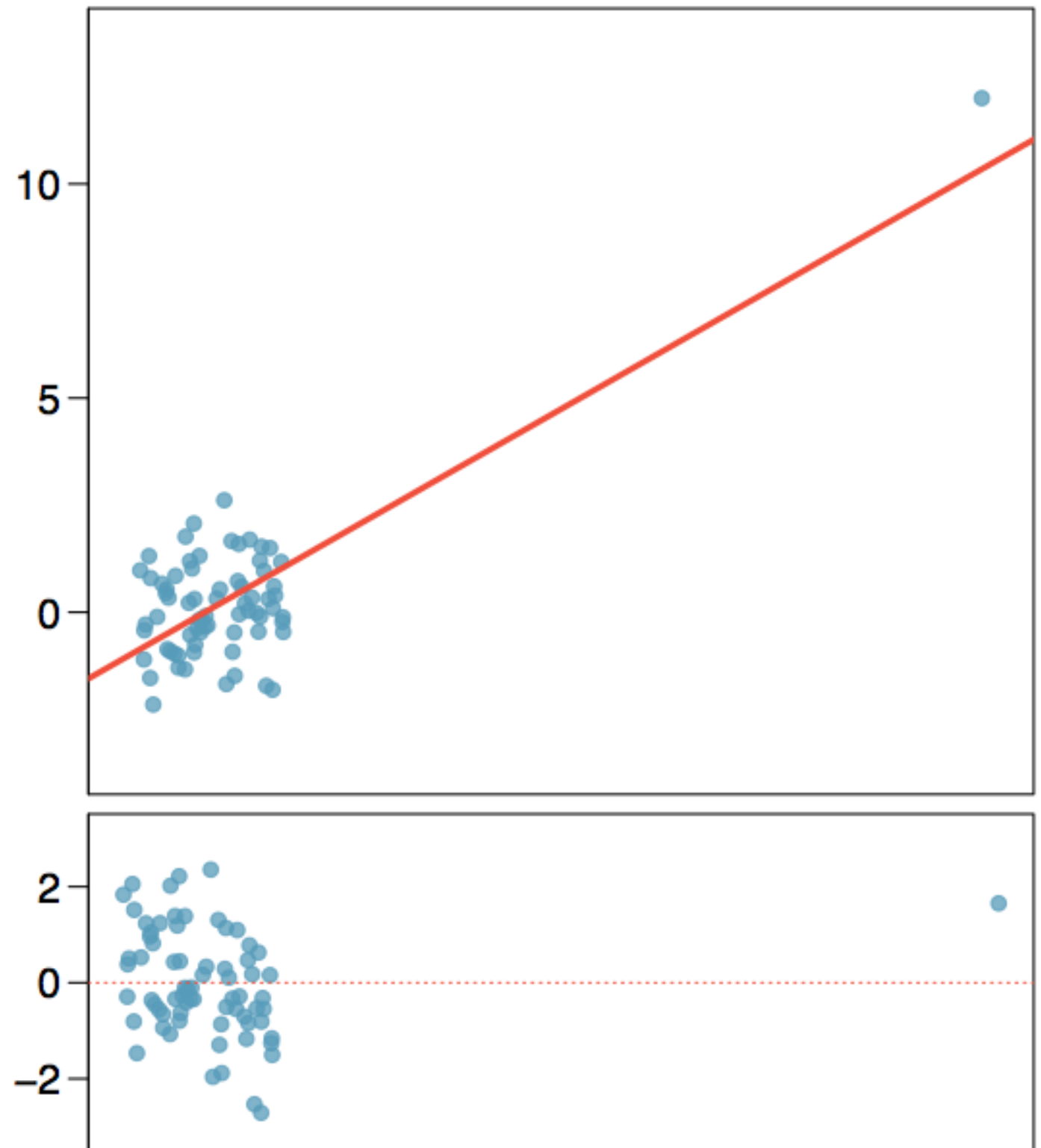
To answer this question think of where the regression line would be with and without the outlier(s).

Without the outliers the regression line would be steeper, and lie closer to the larger group of observations. With the outliers the line is pulled up and away from some of the observations in the larger group.



Types of outliers

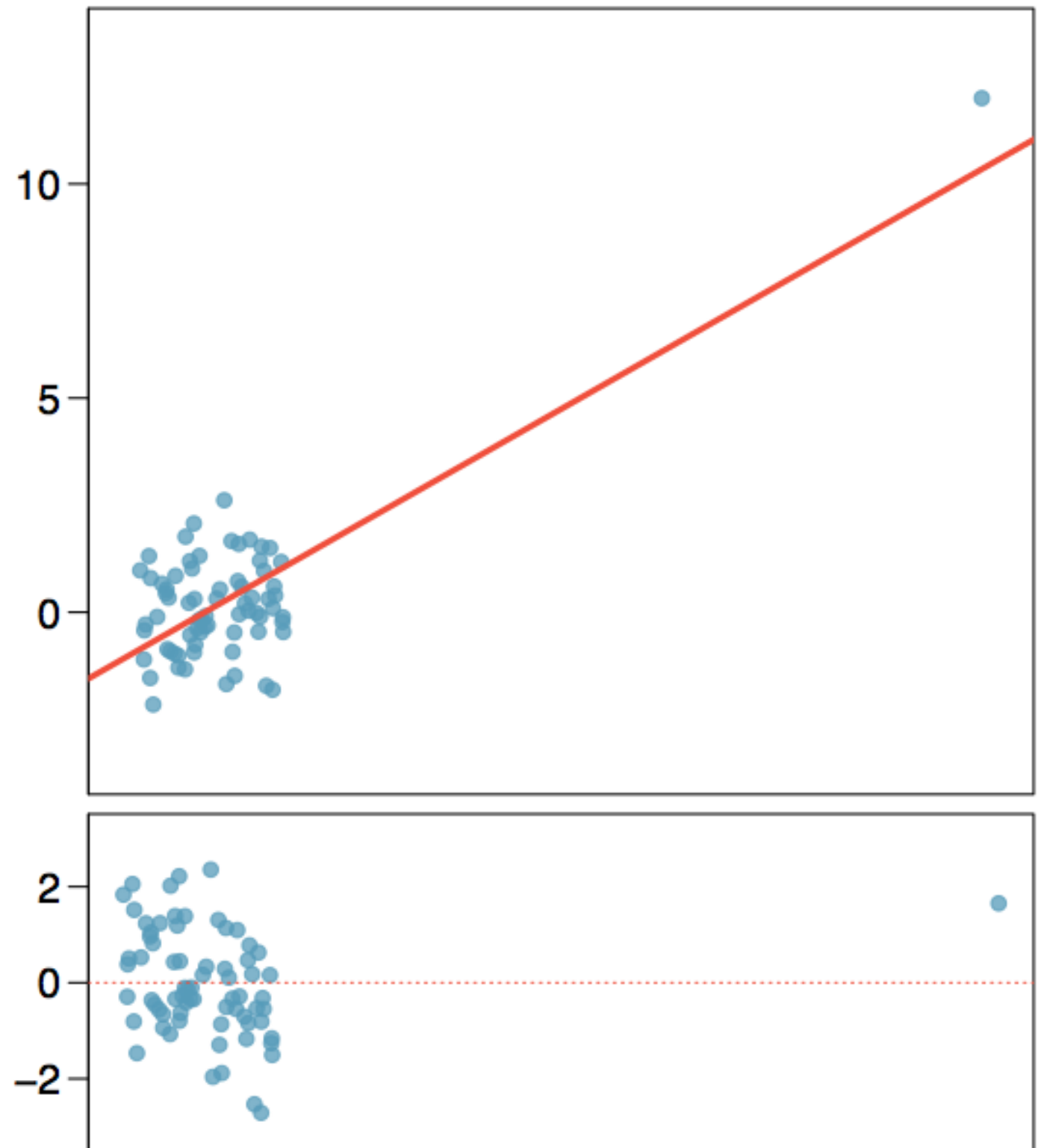
How do outliers influence the least squares line in this plot?



Types of outliers

How do outliers influence the least squares line in this plot?

Without the outlier there is no evident relationship between x and y .



Some terminology

- *Outliers* are points that lie away from the cloud of points.

Some terminology

- *Outliers* are points that lie away from the cloud of points.
- Outliers that lie horizontally away from the center of the cloud are called *high leverage* points.

Some terminology

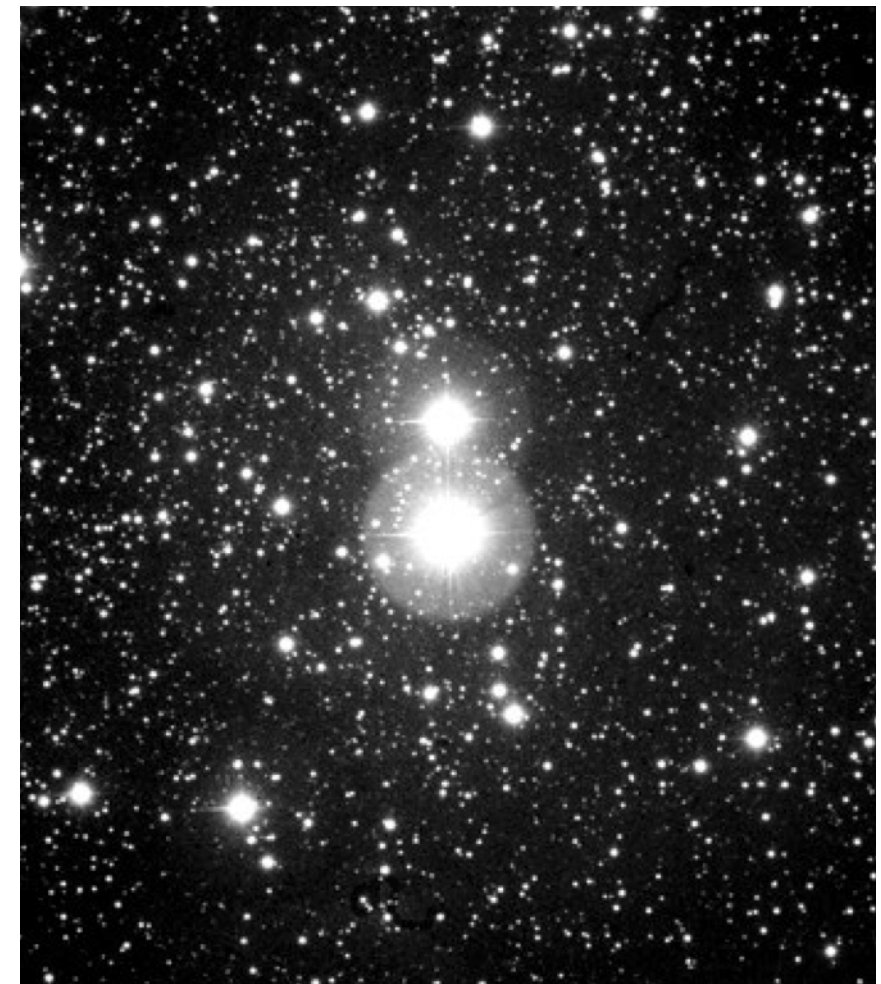
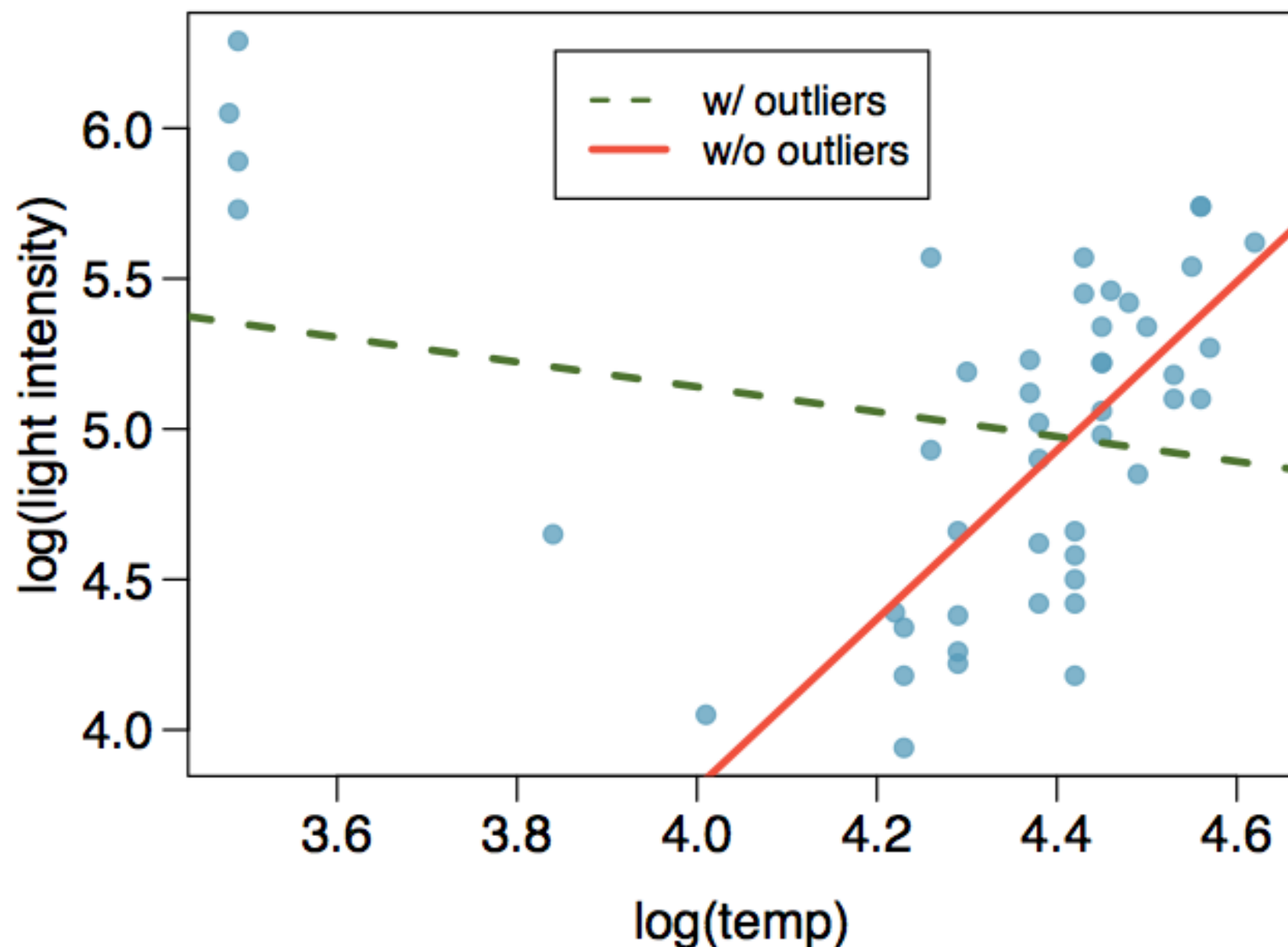
- *Outliers* are points that lie away from the cloud of points.
- Outliers that lie horizontally away from the center of the cloud are called *high leverage* points.
- High leverage points that actually influence the slope of the regression line are called *influential* points.

Some terminology

- *Outliers* are points that lie away from the cloud of points.
- Outliers that lie horizontally away from the center of the cloud are called *high leverage* points.
- High leverage points that actually influence the slope of the regression line are called *influential* points.
- In order to determine if a point is influential, visualize the regression line with and without the point. Does the slope of the line change considerably? If so, then the point is influential. If not, then it's not an influential point.

Influential points

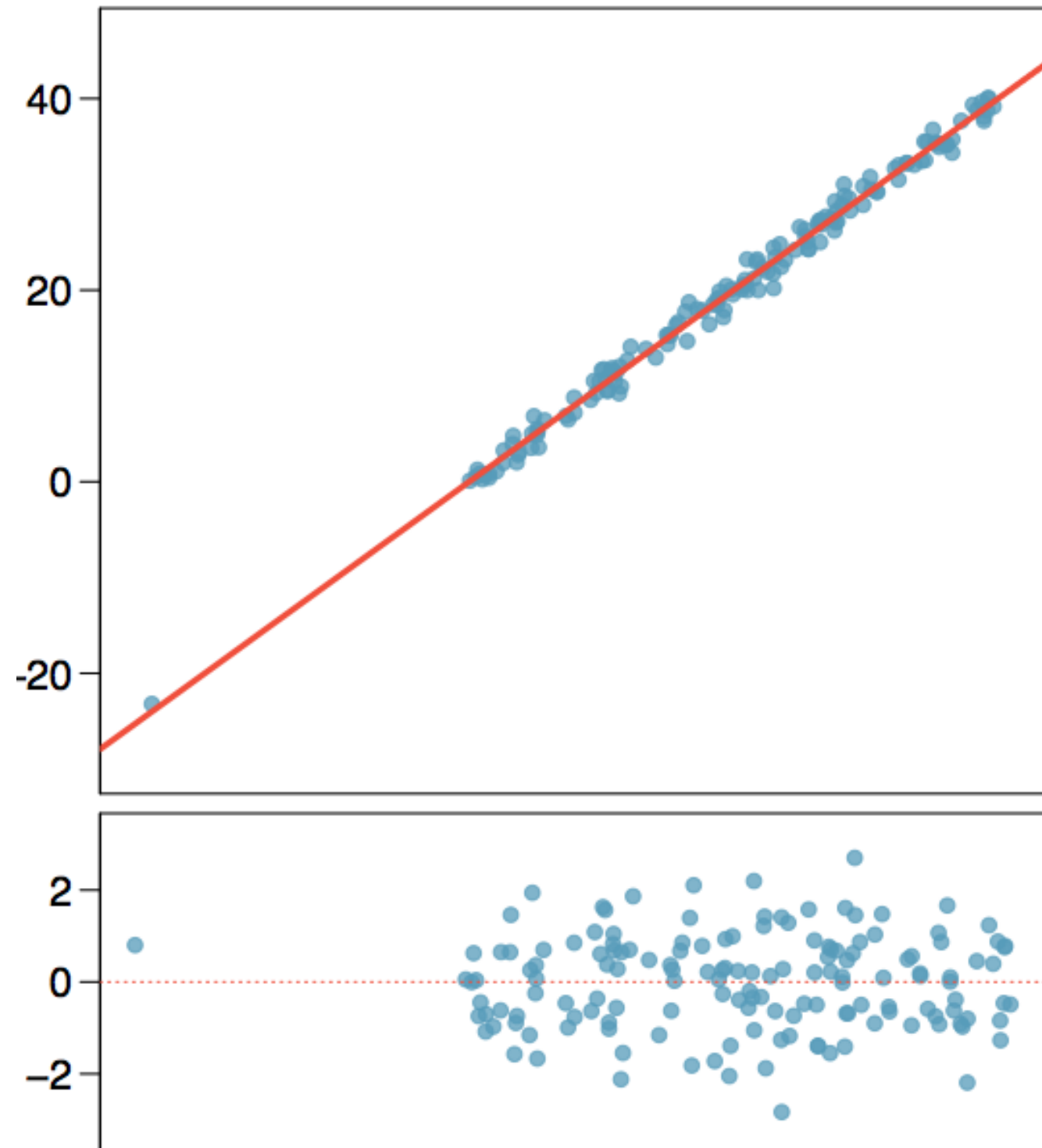
Data are available on the log of the surface temperature and the log of the light intensity of 47 stars in the star cluster CYG OB1.



Types of outliers

Which of the below best describes the outlier?

- (a) influential
- (b) high leverage
- (c) none of the above
- (d) there are no outliers



Types of outliers

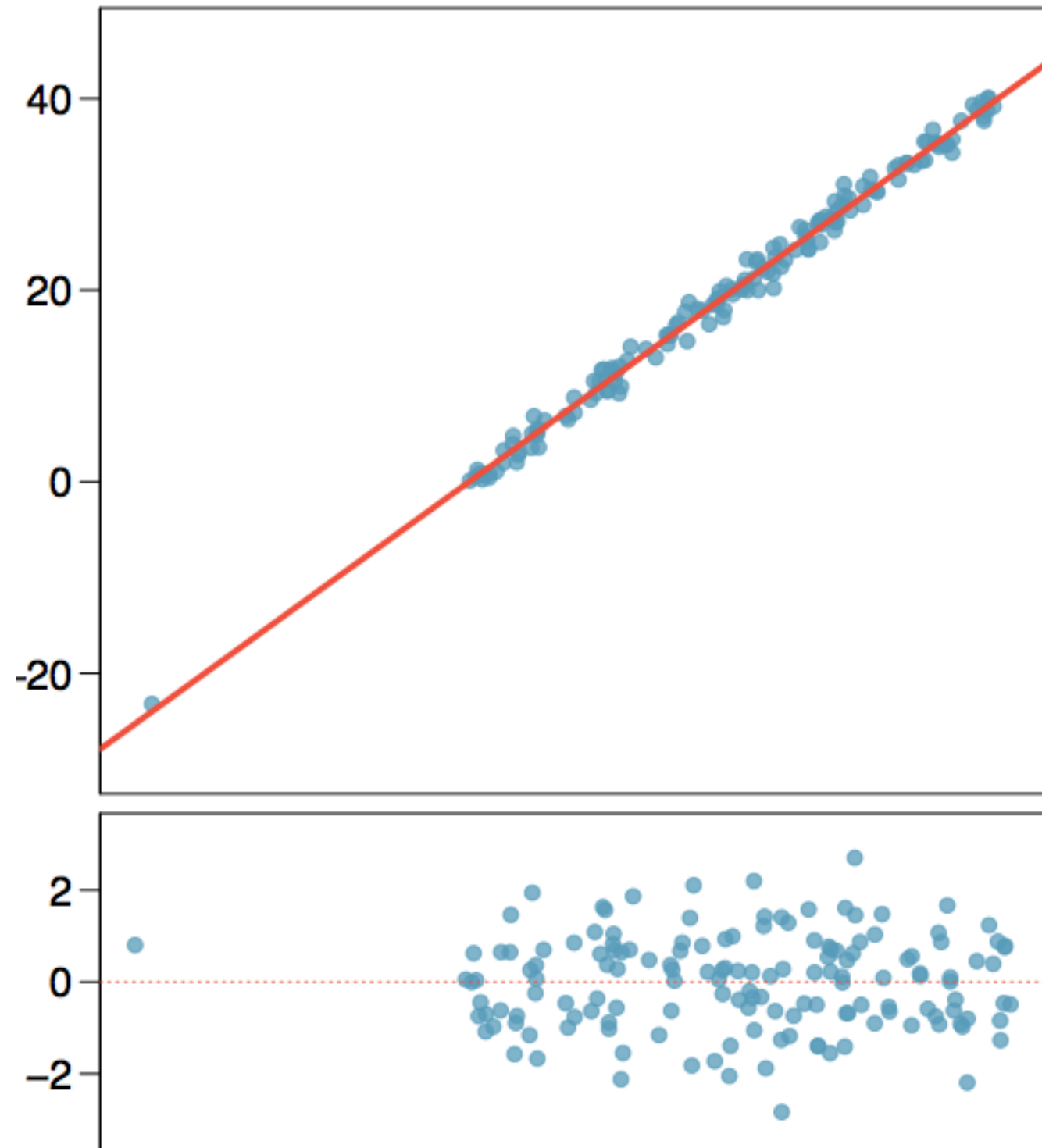
Which of the below best describes the outlier?

(a) influential

(b) high leverage

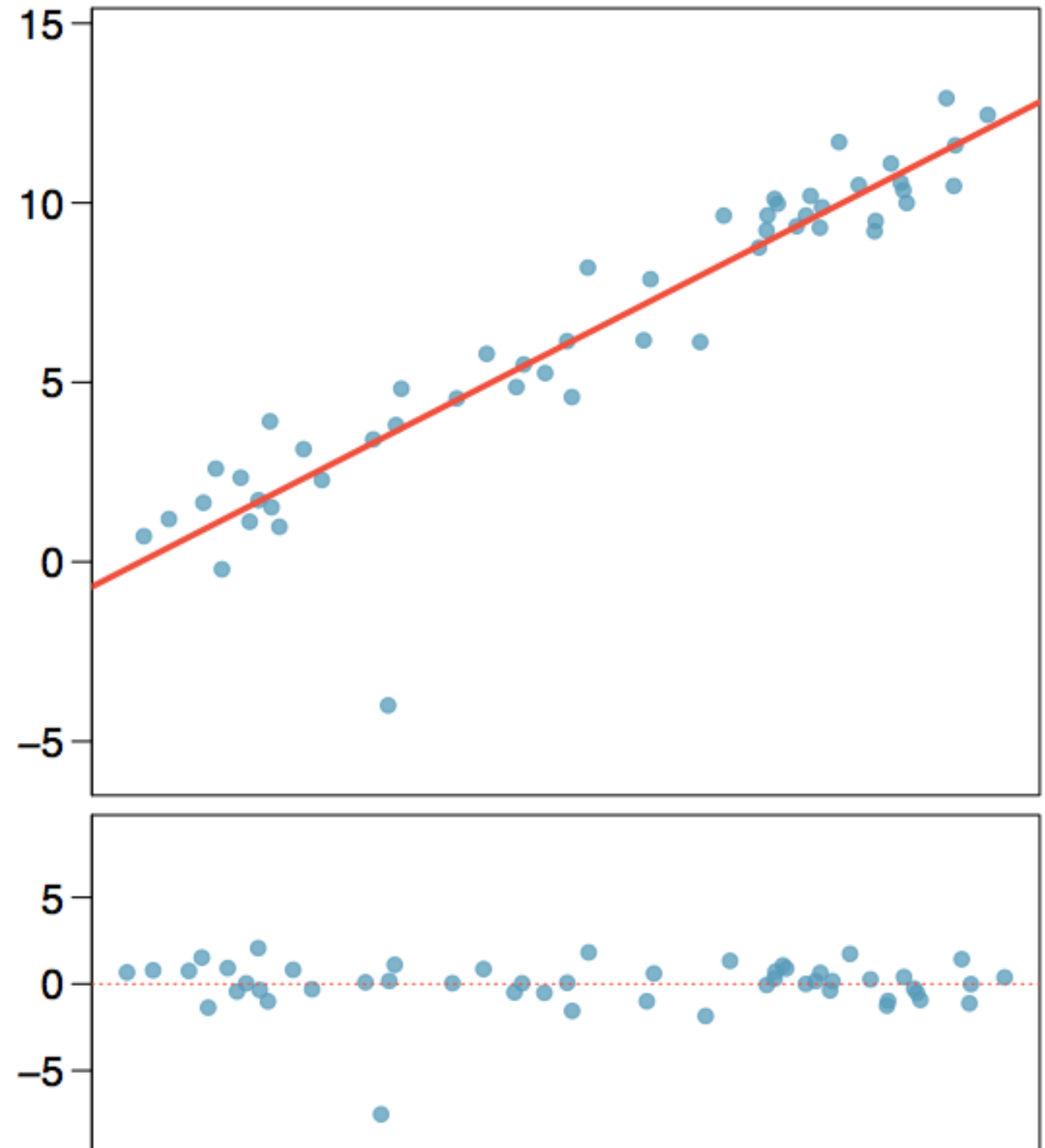
(c) none of the above

(d) there are no outliers



Types of outliers

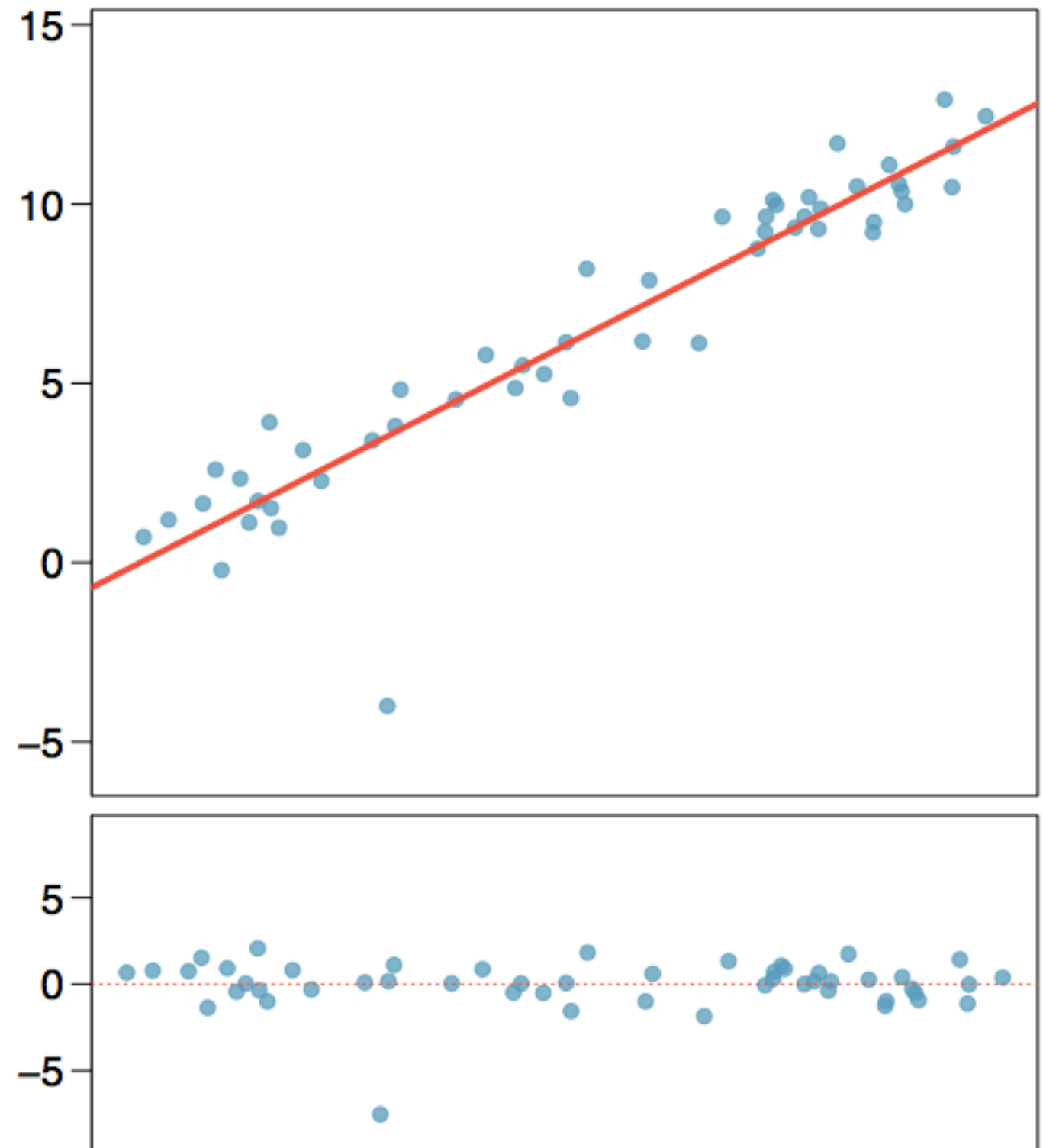
Does this outlier influence the slope of the regression line?



Types of outliers

Does this outlier
influence the slope of
the regression line?

Not much...



Recap

Which of following is true?

- (a) Influential points always change the intercept of the regression line.
- (b) Influential points always reduce R^2 .
- (c) It is much more likely for a low leverage point to be influential, than a high leverage point.
- (d) When the data set includes an influential point, the relationship between the explanatory variable and the response variable is always nonlinear.
- (e) None of the above.

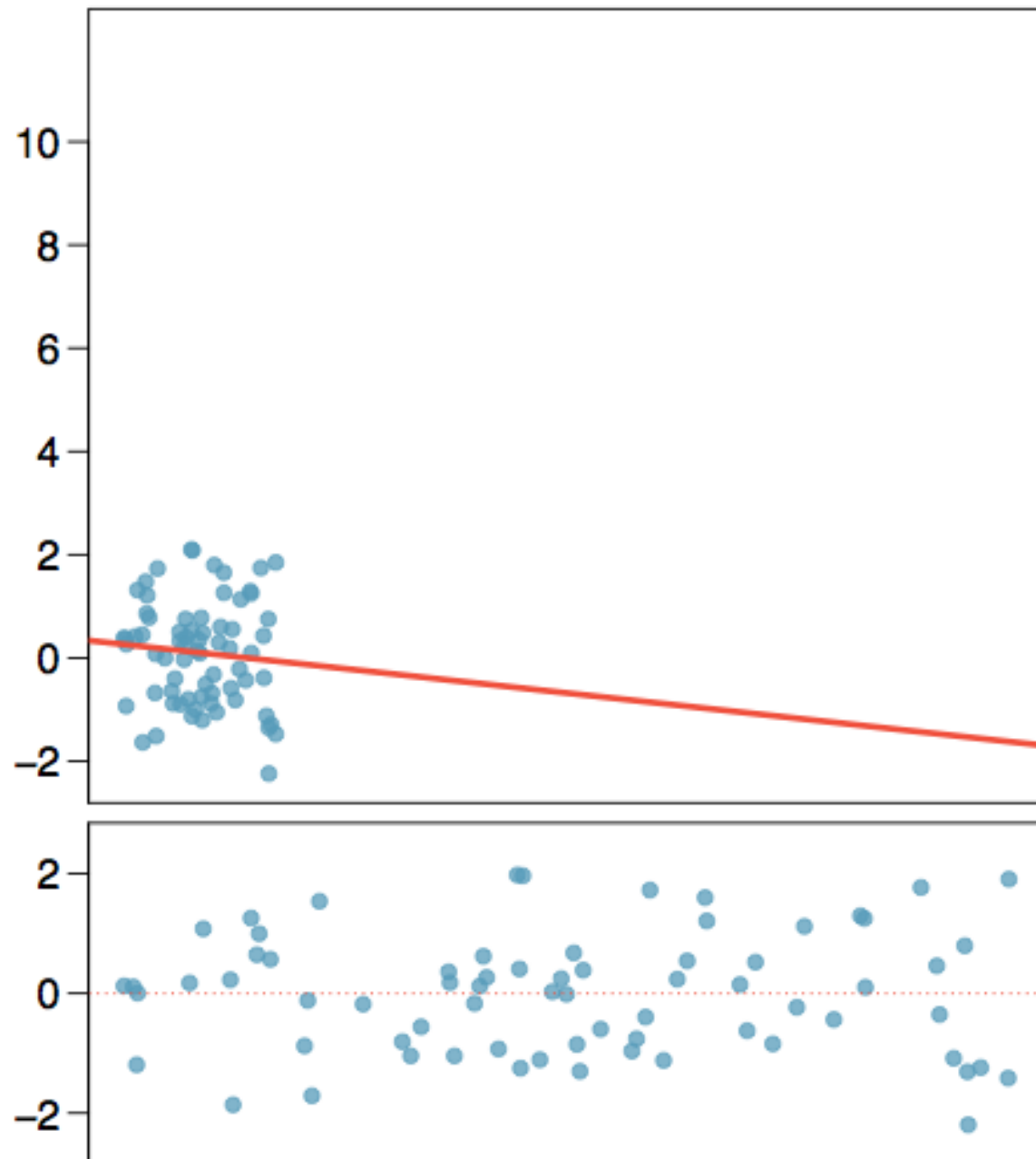
Recap

Which of following is true?

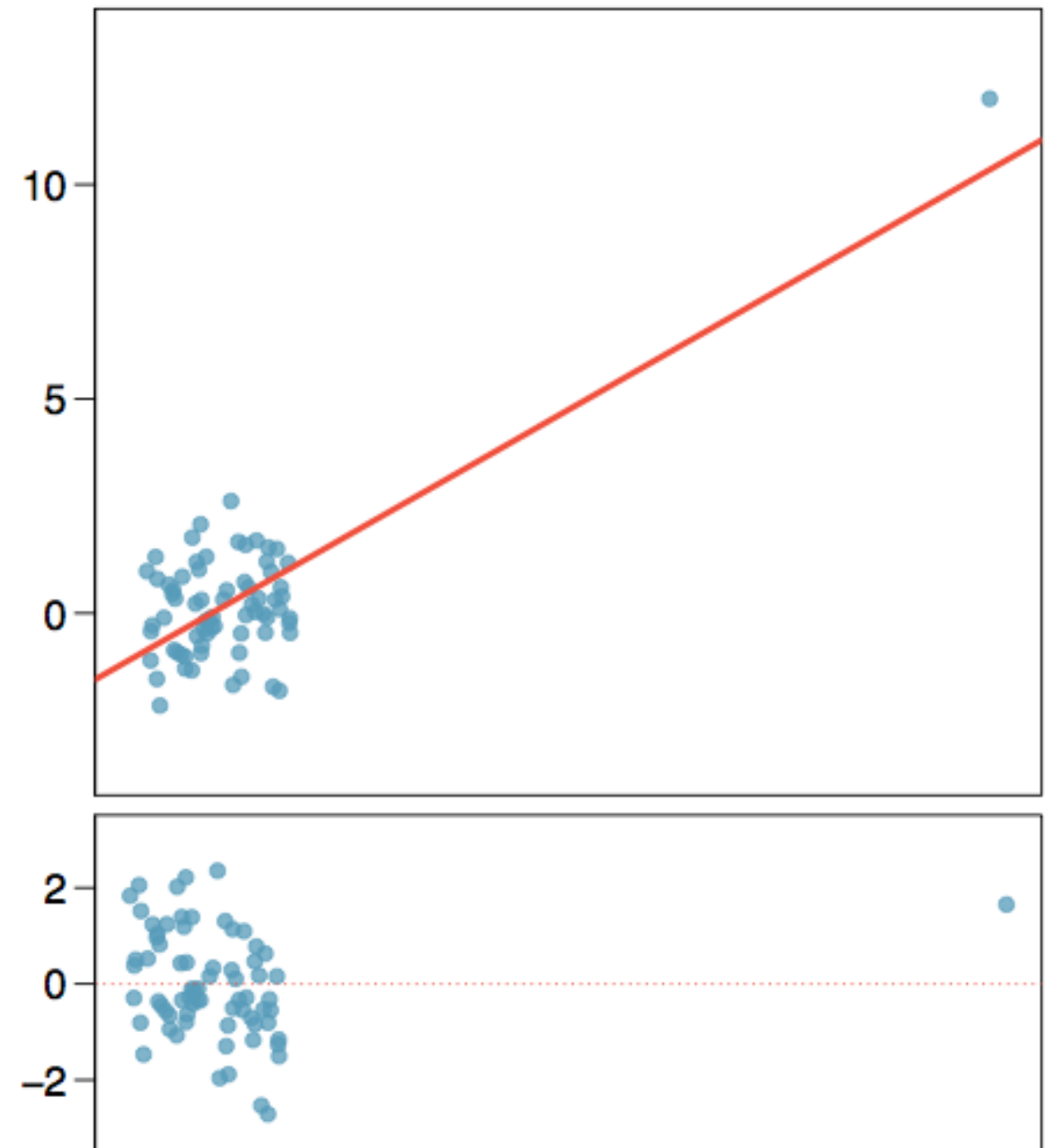
- (a) Influential points always change the intercept of the regression line.
- (b) Influential points always reduce R^2 .
- (c) It is much more likely for a low leverage point to be influential, than a high leverage point.
- (d) When the data set includes an influential point, the relationship between the explanatory variable and the response variable is always nonlinear.
- (e) None of the above.*

Recap (cont.)

$$R = 0.08, R^2 = 0.0064$$



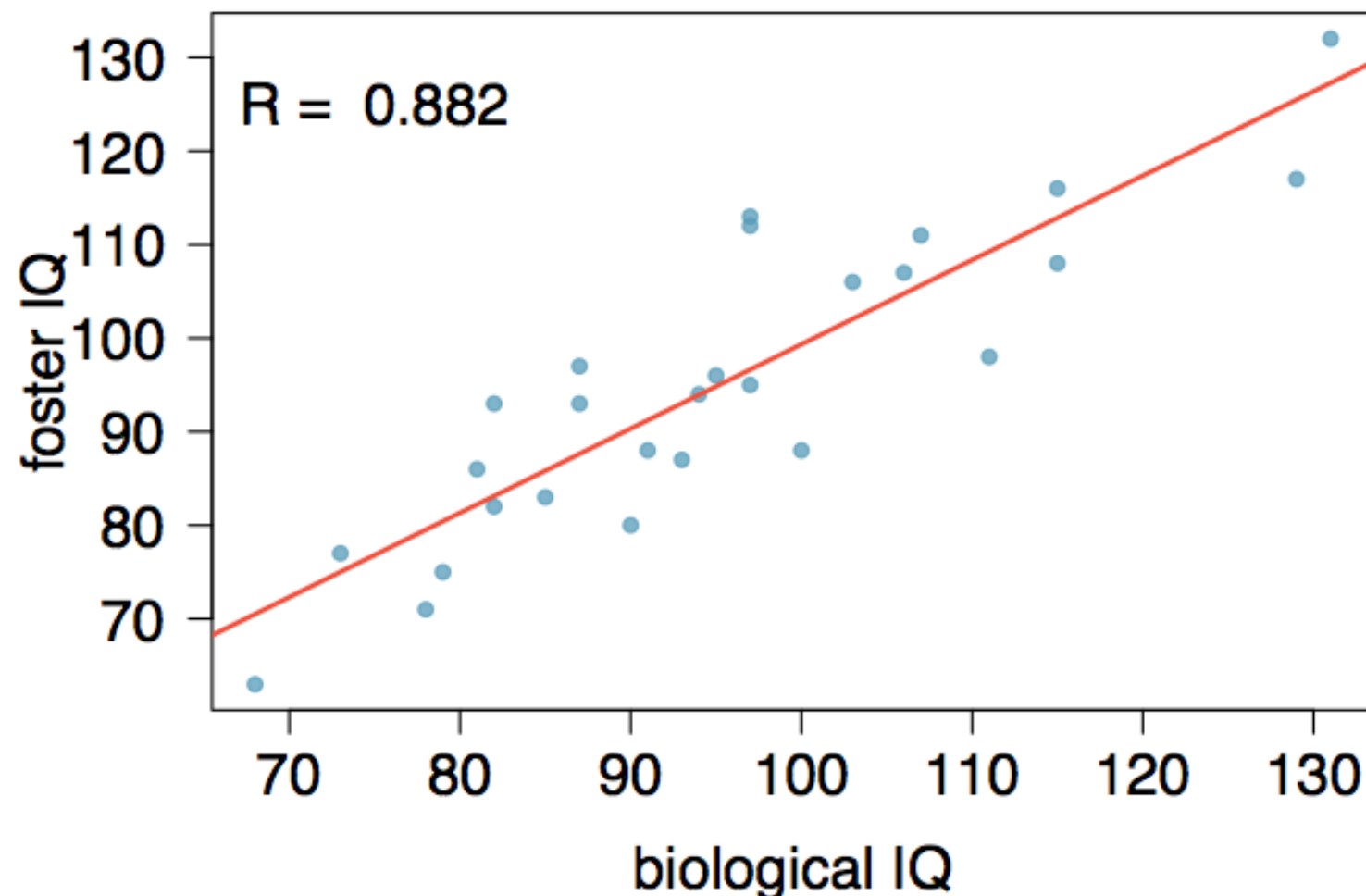
$$R = 0.79, R^2 = 0.6241$$



Inference for Linear Regression

Nature or nurture?

In 1966 Cyril Burt published a paper called "The genetic determination of differences in intelligence: A study of monozygotic twins reared apart?" The data consist of IQ scores for [an assumed random sample of] 27 identical twins, one raised by foster parents, the other by the biological parents.



Practice

Which of the following is false?

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.20760	9.29990	0.990	0.332
bioIQ	0.90144	0.09633	9.358	1.2e-09

Residual standard error: 7.729 on 25 degrees of freedom

Multiple R-squared: 0.7779, Adjusted R-squared: 0.769

F-statistic: 87.56 on 1 and 25 DF, p-value: 1.204e-09

- (a) Additional 10 points in the biological twin's IQ is associated with additional 9 points in the foster twin's IQ, on average.
- (b) Roughly 78% of the foster twins' IQs can be accurately predicted by the model.
- (c) The linear model is $\widehat{fosterIQ} = 9.2 + 0.9 \times bioIQ$
- (d) Foster twins with IQs higher than average IQs tend to have biological twins with higher than average IQs as well.

Practice

Which of the following is false?

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.20760	9.29990	0.990	0.332
bioIQ	0.90144	0.09633	9.358	1.2e-09

Residual standard error: 7.729 on 25 degrees of freedom

Multiple R-squared: 0.7779, Adjusted R-squared: 0.769

F-statistic: 87.56 on 1 and 25 DF, p-value: 1.204e-09

(a) Additional 10 points in the biological twin's IQ is associated with additional 9 points in the foster twin's IQ, on average.

(b) Roughly 78% of the foster twins' IQs can be accurately predicted by the model.

(c) The linear model is $\widehat{fosterIQ} = 9.2 + 0.9 \times bioIQ$

(d) Foster twins with IQs higher than average IQs tend to have biological twins with higher than average IQs as well.

Testing for the slope

Assuming that these 27 twins comprise a representative sample of all twins separated at birth, we would like to test if these data provide convincing evidence that the IQ of the biological twin is a significant predictor of IQ of the foster twin. What are the appropriate hypotheses?

(a) $H_0: b_0 = 0$; $H_A: b_0 \neq 0$

(b) $H_0: \beta_0 = 0$; $H_A: \beta_0 \neq 0$

(c) $H_0: b_1 = 0$; $H_A: b_1 \neq 0$

(d) $H_0: \beta_1 = 0$; $H_A: \beta_1 \neq 0$

Testing for the slope

Assuming that these 27 twins comprise a representative sample of all twins separated at birth, we would like to test if these data provide convincing evidence that the IQ of the biological twin is a significant predictor of IQ of the foster twin. What are the appropriate hypotheses?

(a) $H_0: b_0 = 0$; $H_A: b_0 \neq 0$

(b) $H_0: \beta_0 = 0$; $H_A: \beta_0 \neq 0$

(c) $H_0: b_1 = 0$; $H_A: b_1 \neq 0$

(d) $H_0: \beta_1 = 0$; $H_A: \beta_1 \neq 0$

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

- We always use a t-test in inference for regression.

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

- We always use a t-test in inference for regression.

Remember: test statistic $T = (\text{point estimate} - \text{null value}) / \text{SE}$

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

- We always use a t-test in inference for regression.

Remember: test statistic $T = (\text{point estimate} - \text{null value}) / \text{SE}$

- Point estimate = b_1 is the observed slope.

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

- We always use a t-test in inference for regression.

Remember: test statistic $T = (\text{point estimate} - \text{null value}) / \text{SE}$

- Point estimate = b_1 is the observed slope.
- SE_{b_1} is the standard error associated with the slope.

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

- We always use a t-test in inference for regression.

Remember: test statistic $T = (\text{point estimate} - \text{null value}) / \text{SE}$

- Point estimate = b_1 is the observed slope.
- SE_{b_1} is the standard error associated with the slope.
- Degrees of freedom associated with the slope is $df = n - 2$, where n is the sample size.

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

- We always use a t-test in inference for regression.

Remember: test statistic $T = (\text{point estimate} - \text{null value}) / \text{SE}$

- Point estimate = b_1 is the observed slope.
- SE_{b_1} is the standard error associated with the slope.
- Degrees of freedom associated with the slope is $df = n - 2$, where n is the sample size.

Remember: we lose 1 degree of freedom for each parameter we estimate, and in simple linear regression we estimate 2 parameters, β_0 and β_1 .

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

$$T = \frac{0.9014 - 0}{0.0963} = 9.36$$

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

$$T = \frac{0.9014 - 0}{0.0963} = 9.36$$

$$df = 27 - 2 = 25$$

Testing for the slope (cont.)

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
bioIQ	0.9014	0.0963	9.36	0.0000

$$T = \frac{0.9014 - 0}{0.0963} = 9.36$$

$$df = 27 - 2 = 25$$

$$p - value = P(|T| > 9.36) < 0.01$$

Confidence interval for the slope

Remember that a confidence interval is calculated as *point estimate* \pm *ME* and the degrees of freedom associated with the slope in a simple linear regression is $n - 2$. Which of the below is the correct 95% confidence interval for the slope parameter? Note that the model is based on observations from 27 twins.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
biolQ	0.9014	0.0963	9.36	0.0000

- (a) $9.2076 \pm 1.65 \times 9.2999$
- (b) $0.9014 \pm 2.06 \times 0.0963$
- (c) $0.9014 \pm 1.96 \times 0.0963$
- (d) $9.2076 \pm 1.96 \times 0.0963$

Confidence interval for the slope

Remember that a confidence interval is calculated as *point estimate* \pm *ME* and the degrees of freedom associated with the slope in a simple linear regression is $n - 2$. Which of the below is the correct 95% confidence interval for the slope parameter? Note that the model is based on observations from 27 twins.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
biolQ	0.9014	0.0963	9.36	0.0000

$$n = 27 \quad df = 27 - 2 = 25$$

- (a) $9.2076 \pm 1.65 \times 9.2999$
- (b) $0.9014 \pm 2.06 \times 0.0963$
- (c) $0.9014 \pm 1.96 \times 0.0963$
- (d) $9.2076 \pm 1.96 \times 0.0963$

Confidence interval for the slope

Remember that a confidence interval is calculated as *point estimate* \pm *ME* and the degrees of freedom associated with the slope in a simple linear regression is $n - 2$. Which of the below is the correct 95% confidence interval for the slope parameter? Note that the model is based on observations from 27 twins.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
biolQ	0.9014	0.0963	9.36	0.0000

$$n = 27 \quad df = 27 - 2 = 25$$

$$95\%: t_{25}^* = 2.06$$

(a) $9.2076 \pm 1.65 \times 9.2999$

(b) $0.9014 \pm 2.06 \times 0.0963$

(c) $0.9014 \pm 1.96 \times 0.0963$

(d) $9.2076 \pm 1.96 \times 0.0963$

Confidence interval for the slope

Remember that a confidence interval is calculated as *point estimate* \pm *ME* and the degrees of freedom associated with the slope in a simple linear regression is $n - 2$. Which of the below is the correct 95% confidence interval for the slope parameter? Note that the model is based on observations from 27 twins.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
biolQ	0.9014	0.0963	9.36	0.0000

(a) $9.2076 \pm 1.65 \times 9.2999$

(b) $0.9014 \pm 2.06 \times 0.0963$

(c) $0.9014 \pm 1.96 \times 0.0963$

(d) $9.2076 \pm 1.96 \times 0.0963$

$$n = 27 \quad df = 27 - 2 = 25$$

$$95\%: t_{25}^* = 2.06$$

$$0.9014 \pm 2.06 \times 0.0963$$

Confidence interval for the slope

Remember that a confidence interval is calculated as *point estimate* \pm *ME* and the degrees of freedom associated with the slope in a simple linear regression is $n - 2$. Which of the below is the correct 95% confidence interval for the slope parameter? Note that the model is based on observations from 27 twins.

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.2076	9.2999	0.99	0.3316
biolQ	0.9014	0.0963	9.36	0.0000

(a) $9.2076 \pm 1.65 \times 9.2999$

(b) $0.9014 \pm 2.06 \times 0.0963$

(c) $0.9014 \pm 1.96 \times 0.0963$

(d) $9.2076 \pm 1.96 \times 0.0963$

$$n = 27 \quad df = 27 - 2 = 25$$

$$95\%: t_{25}^* = 2.06$$

$$0.9014 \pm 2.06 \times 0.0963$$

$$(0.7, 1.1)$$

Recap

- Inference for the slope for a single-predictor linear regression model:

Recap

- Inference for the slope for a single-predictor linear regression model:
- Hypothesis test:

$$T = \frac{b_1 - \text{null value}}{SE_{b_1}} \quad df = n - 2$$

Recap

- Inference for the slope for a single-predictor linear regression model:

- Hypothesis test:

$$T = \frac{b_1 - \text{null value}}{SE_{b_1}} \quad df = n - 2$$

- Confidence interval:

$$b_1 \pm t_{df=n-2}^* SE_{b_1}$$

Recap

- Inference for the slope for a single-predictor linear regression model:

- Hypothesis test:

$$T = \frac{b_1 - \text{null value}}{SE_{b_1}} \quad df = n - 2$$

- Confidence interval:

$$b_1 \pm t_{df=n-2}^* SE_{b_1}$$

- The null value is often 0 since we are usually checking for *any* relationship between the explanatory and the response variable.

Recap

- Inference for the slope for a single-predictor linear regression model:

- Hypothesis test:

$$T = \frac{b_1 - \text{null value}}{SE_{b_1}} \quad df = n - 2$$

- Confidence interval:

$$b_1 \pm t_{df=n-2}^* SE_{b_1}$$

- The null value is often 0 since we are usually checking for *any* relationship between the explanatory and the response variable.
- The regression output gives b_1 , SE_{b_1} , and *two-tailed* p-value for the t-test for the slope where the null value is 0.

Recap

- Inference for the slope for a single-predictor linear regression model:

- Hypothesis test:

$$T = \frac{b_1 - \text{null value}}{SE_{b_1}} \quad df = n - 2$$

- Confidence interval:

$$b_1 \pm t_{df=n-2}^* SE_{b_1}$$

- The null value is often 0 since we are usually checking for *any* relationship between the explanatory and the response variable.
- The regression output gives b_1 , SE_{b_1} , and *two-tailed* p-value for the t-test for the slope where the null value is 0.
- We rarely do inference on the intercept, so we'll be focusing on the estimates and inference for the slope.

Caution

- Always be aware of the type of data you're working with: random sample, non-random sample, or population.

Caution

- Always be aware of the type of data you're working with: random sample, non-random sample, or population.
- Statistical inference, and the resulting p-values, are meaningless when you already have population data.

Caution

- Always be aware of the type of data you're working with: random sample, non-random sample, or population.
- Statistical inference, and the resulting p-values, are meaningless when you already have population data.
- If you have a sample that is non-random (biased), inference on the results will be unreliable.

Caution

- Always be aware of the type of data you're working with: random sample, non-random sample, or population.
- Statistical inference, and the resulting p-values, are meaningless when you already have population data.
- If you have a sample that is non-random (biased), inference on the results will be unreliable.
- The ultimate goal is to have independent observations.

Thursday is R Session!

- The first half of the lecture will be R session
- The second half of the lecture will be the time for the project discussion
- Objective: Hypothesis Testing
 - Try to formally test your question/ hypotheses from exploratory data analysis(EDA)
 - Interpret your conclusion and present your finding in an appealing way