

CHAPTER 6: INTRODUCTION TO LINEAR REGRESSION

math 189 : data analysis and inference : winter 2018

Jelena Bradic

<http://www.math.ucsd.edu/~jbradic/>

Assistant Professor, Department of Mathematics, University of California, San Diego

jbradic@ucsd.edu

Line fitting, residuals, and correlation

Fitting a line by least squares regression

Types of outliers in linear regression

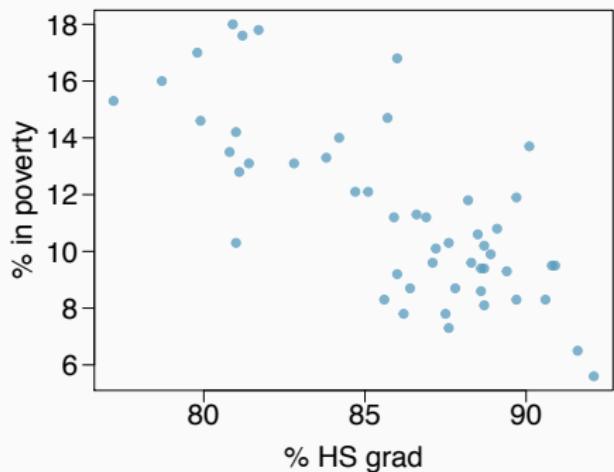
Inference for linear regression

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 % 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

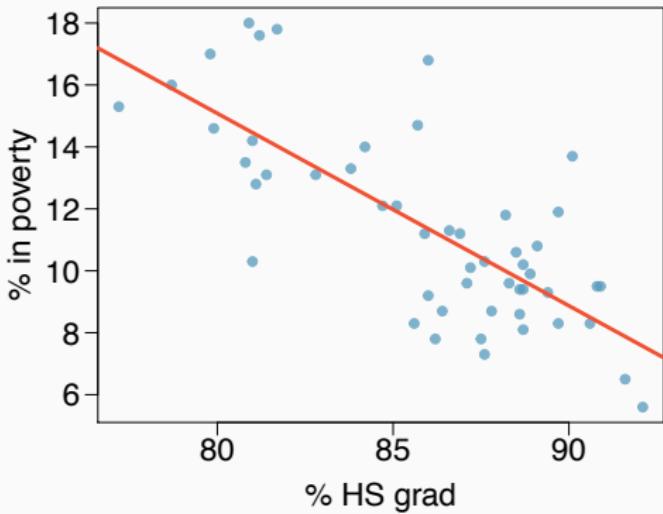
QUANTIFYING THE RELATIONSHIP

- * 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 % in poverty and % HS grad?

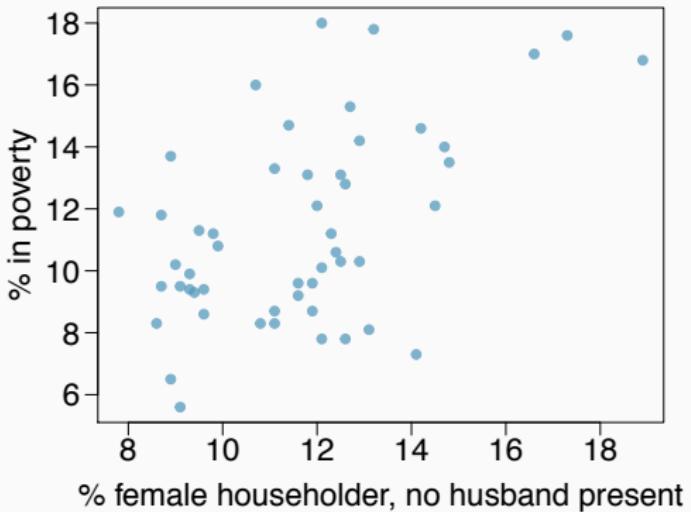
- * 0.6
- * **-0.75**
- * -0.1
- * 0.02
- * -1.5



GUESSING THE CORRELATION

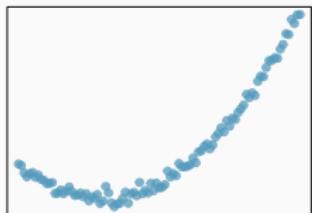
Which of the following is the best guess for the correlation between % in poverty and % HS grad?

- * 0.1
- * -0.6
- * -0.4
- * 0.9
- * 0.5

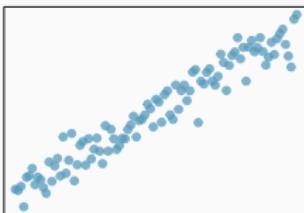


ASSESSING THE CORRELATION

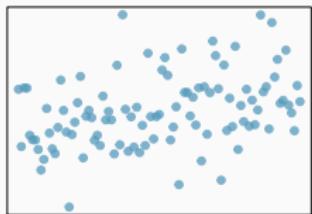
Which of the following has the strongest correlation, i.e. correlation coefficient closest to +1 or -1?



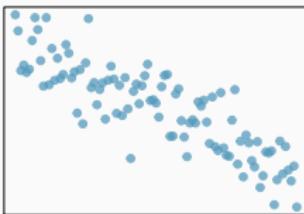
(a)



(b)



(c)



(d)

(b) →
correlation
means linear
association

Line fitting, residuals, and correlation

Fitting a line by least squares regression

Eyeballing the line

Residuals

Best line

The least squares line

Recap: Interpreting the slope and the intercept

Prediction & extrapolation

Conditions for the least squares line

R^2

Categorical explanatory variables

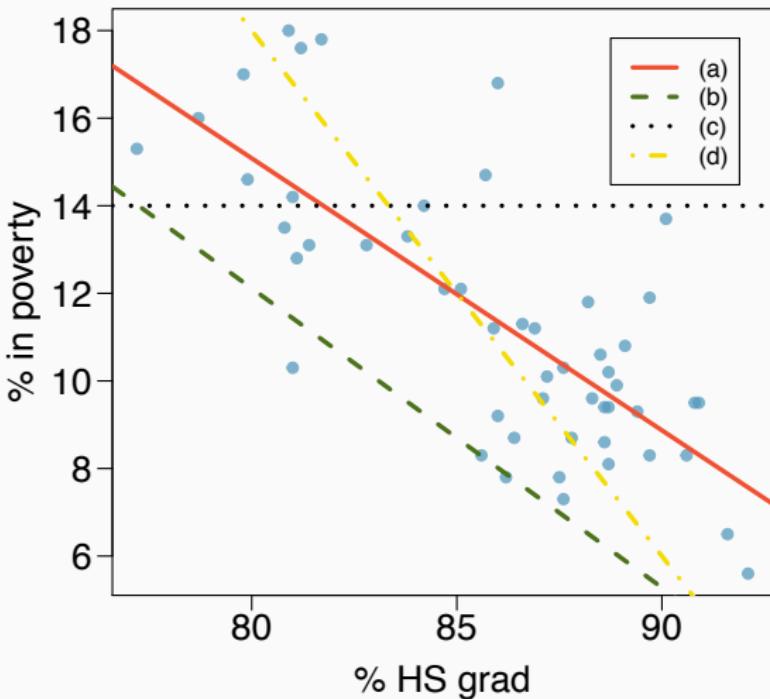
Types of outliers in linear regression

Inference for linear regression

EYEBALLING THE LINE

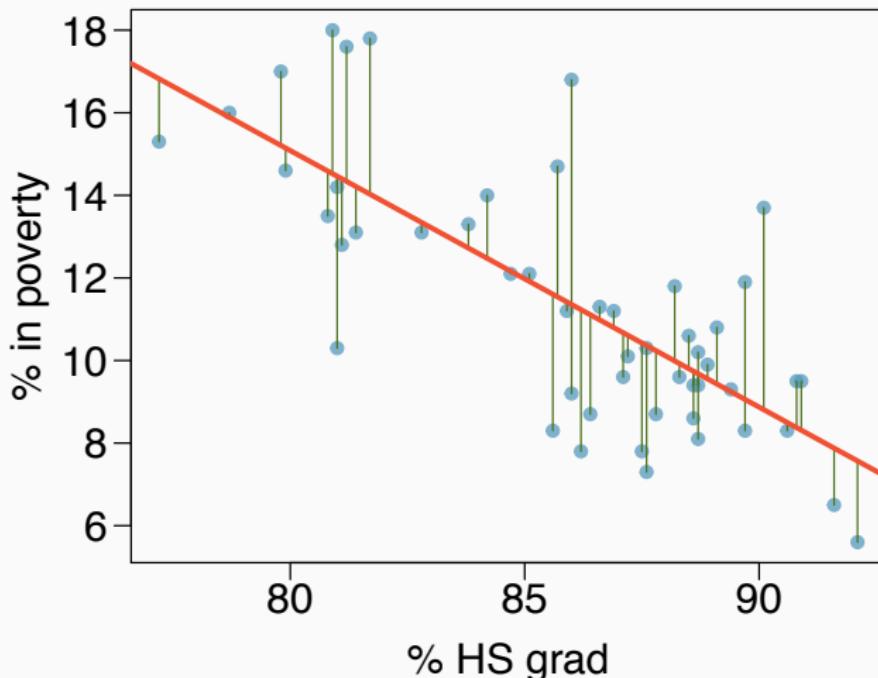
Which of the following appears to be the line that best fits the linear relationship between % in poverty and % HS grad? Choose one.

(a)



RESIDUALS

Residuals are the leftovers from the model fit: Data = Fit + Residual

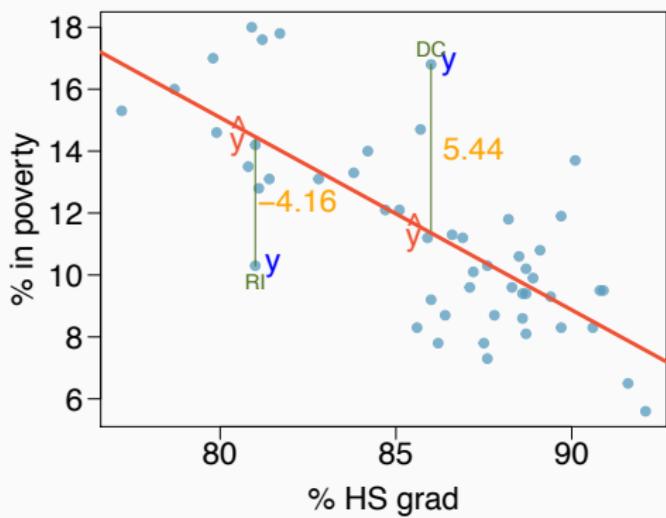


RESIDUALS (CONT.)

Residual

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.

A MEASURE FOR THE BEST LINE

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

$$|e_1| + |e_2| + \cdots + |e_n|$$

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

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

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

THE LEAST SQUARES LINE

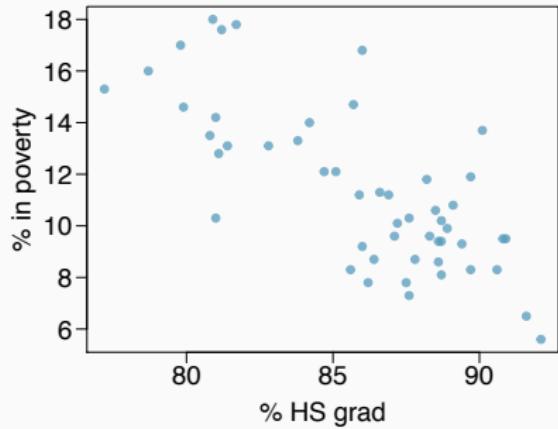
$$\hat{y} = \beta_0 + \beta_1 x$$

The diagram illustrates the components of the linear regression equation $\hat{y} = \beta_0 + \beta_1 x$. It features a central horizontal line labeled "predicted y". Two arrows point from the left towards the top-left corner of the line, labeled "intercept" below the line. Another arrow points from the right towards the middle of the line, labeled "slope" to its left and "explanatory variable" to its right.

Notation:

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

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

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

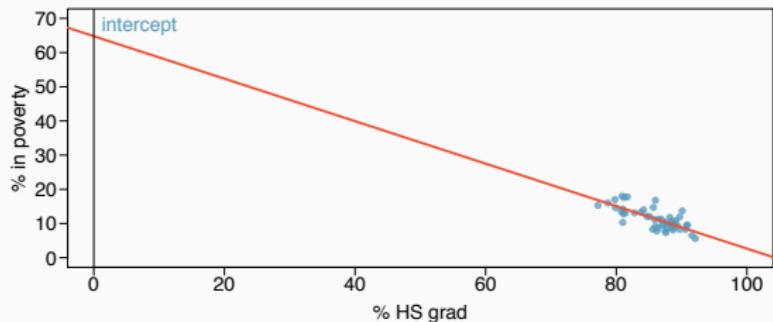
For each additional % point in HS graduate rate, we would expect the % living in poverty to be lower on average by 0.62% points.

INTERCEPT

Intercept

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

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



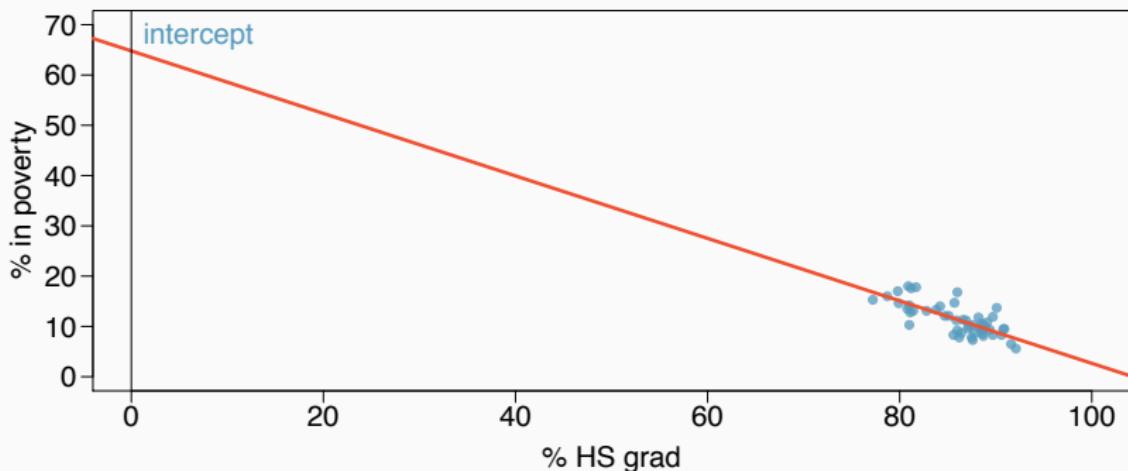
$$\begin{aligned} b_0 &= 11.35 - (-0.62) \times 86.01 \\ &= 64.68 \end{aligned}$$

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

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

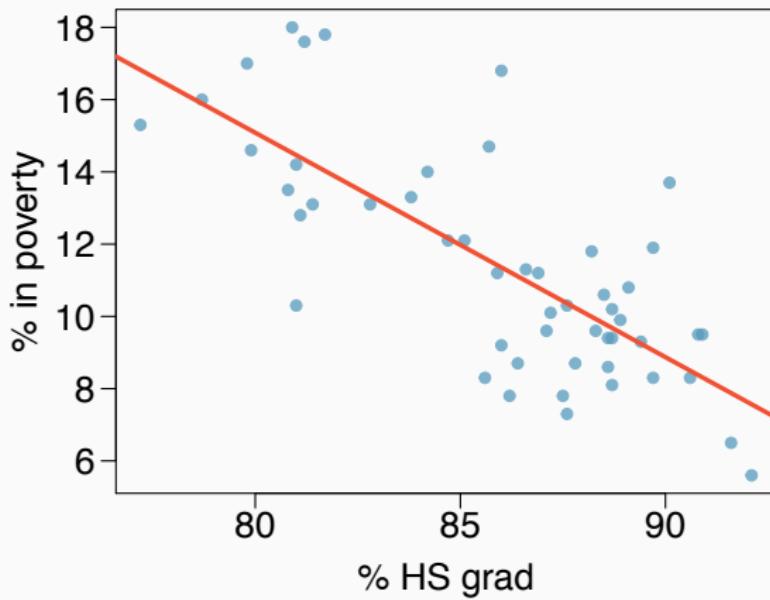
MORE ON THE INTERCEPT

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.



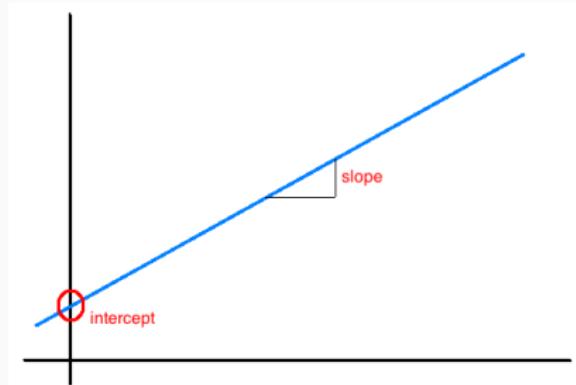
REGRESSION LINE

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



INTERPRETATION OF SLOPE AND INTERCEPT

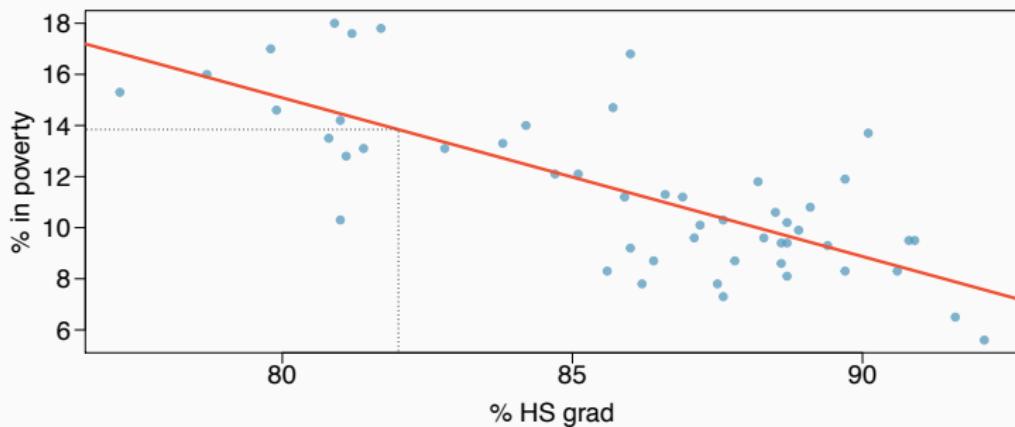
- * **Intercept:** When $x = 0$, y is expected to equal the intercept.
- * **Slope:** For each unit in x , y is expected to increase / decrease on average by the slope.



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

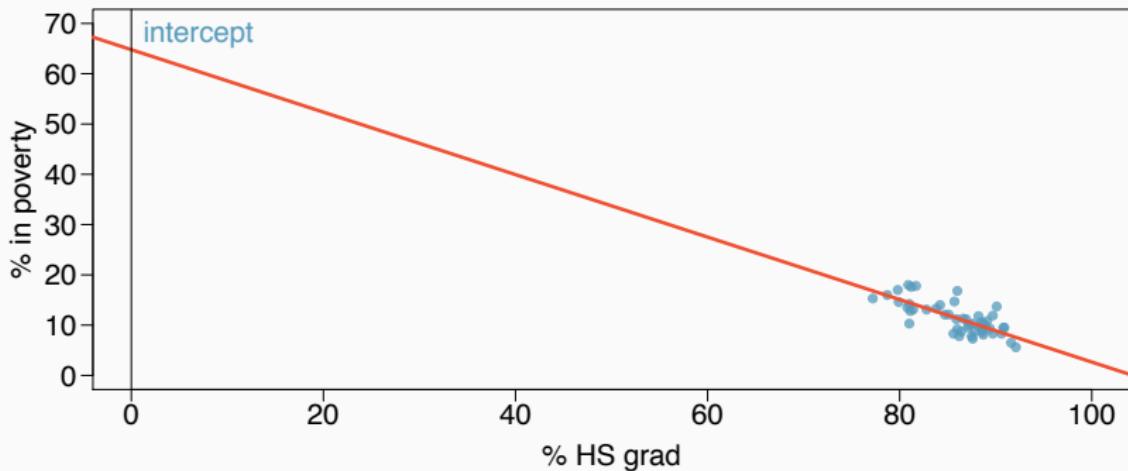
PREDICTION

- * 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.

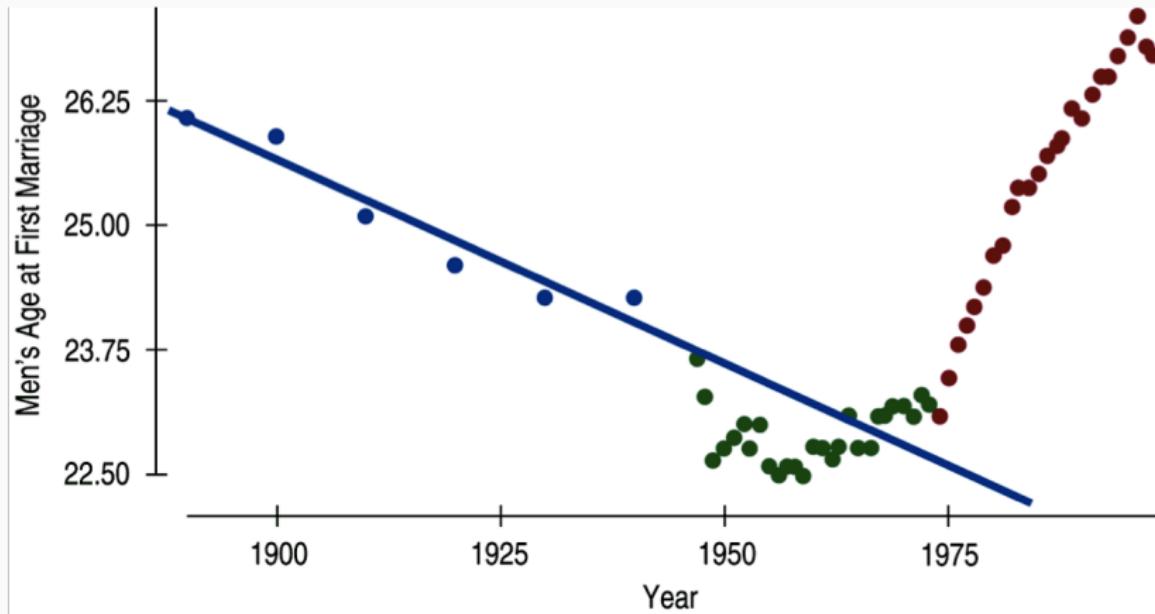


EXTRAPOLATION

- * 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



EXAMPLES OF EXTRAPOLATION

BBC
NEWS

[News Front Page](#)

Last Updated: Thursday, 30 September, 2004, 04:04 GMT 05:04 UK

[E-mail this to a friend](#)

[Printable version](#)

Women 'may outsprint men by 2156'

Women sprinters may be outrunning men in the 2156 Olympics if they continue to close the gap at the rate they are doing, according to scientists.

An Oxford University study found that women are running faster than they have ever done over 100m.

At their current rate of improvement, they should overtake men within 150 years, said Dr Andrew Tatem.

The study, comparing winning times for the Olympic 100m since 1900, is published in the journal Nature.

However, former British Olympic sprinter Derek Redmond told the BBC: "I find it difficult to believe.

"I can see the gap closing between men and women but I can't necessarily see it being overtaken because mens' times are also going to improve."



Women are set to become the dominant sprinters

Momentous sprint at the 2156 Olympics?

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

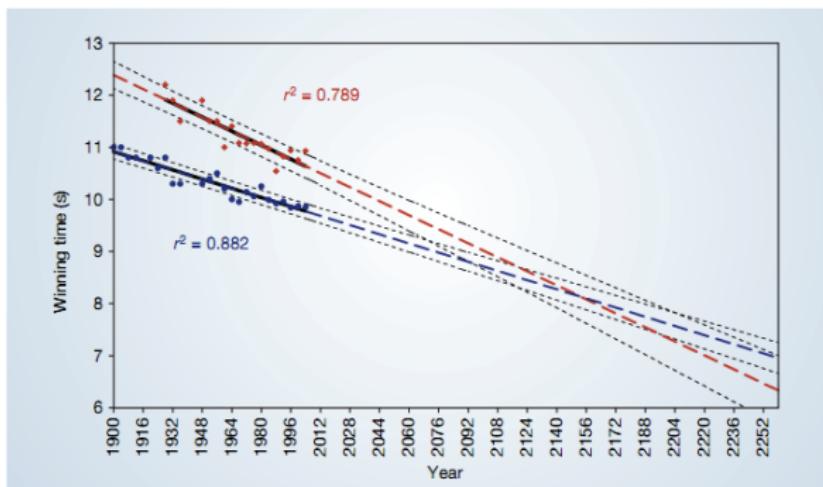


Figure 1 The winning Olympic 100-metre sprint times for men (blue points) and women (red points), with superimposed best-fit linear regression lines (solid black lines) and coefficients of determination. The regression lines are extrapolated (broken blue and red lines for men and women, respectively) 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.

CONDITIONS FOR THE LEAST SQUARES LINE

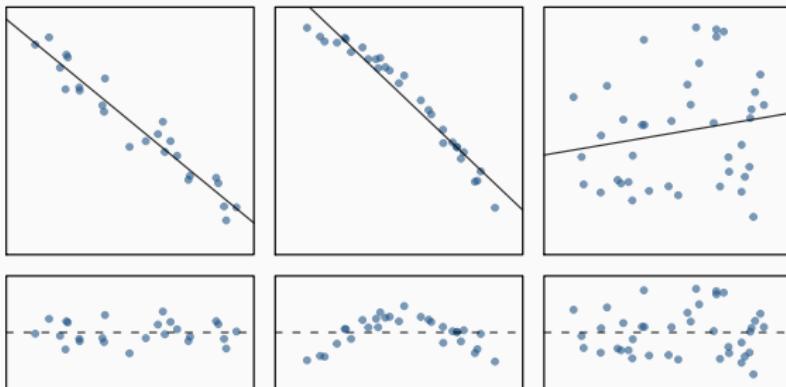
- * Linearity
- * Nearly normal residuals
- * Constant variability

CONDITIONS: (1) LINEARITY

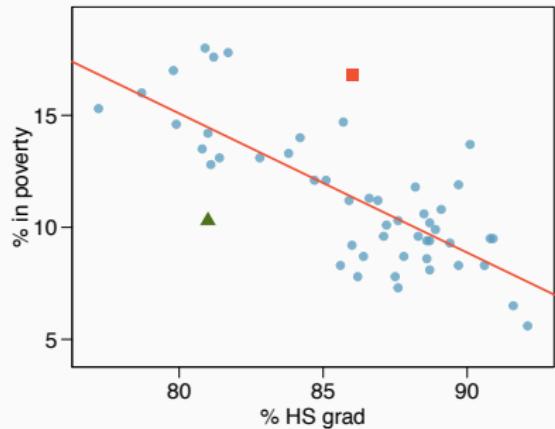
- * The relationship between the explanatory and the response variable should be linear.
- * Methods for fitting a model to non-linear relationships exist, but are beyond the scope of this class. If this topic is of interest, an Online Extra is available on openintro.org covering new techniques.

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. If this topic is of interest, an Online Extra is available on openintro.org covering new techniques.
- * Check using a scatterplot of the data, or a [residuals plot](#).



ANATOMY OF A RESIDUALS PLOT

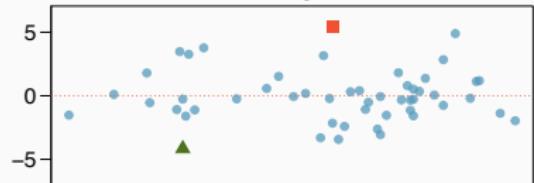


▲ RI:

$$\% \text{ HS grad} = 81 \quad \% \text{ in poverty} = 10.3$$

$$\widehat{\% \text{ in poverty}} = 64.68 - 0.62 * 81 = 14.46$$

$$\begin{aligned} e &= \% \text{ in poverty} - \widehat{\% \text{ in poverty}} \\ &= 10.3 - 14.46 = -4.16 \end{aligned}$$



■ DC:

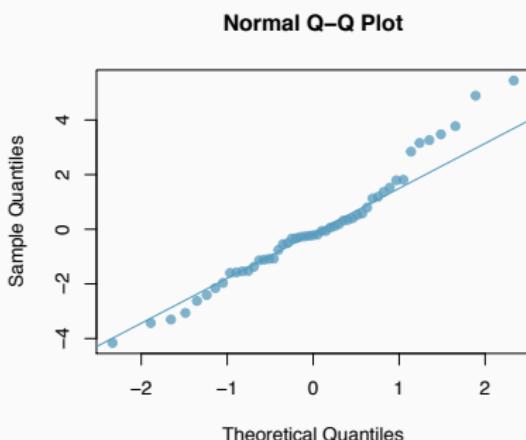
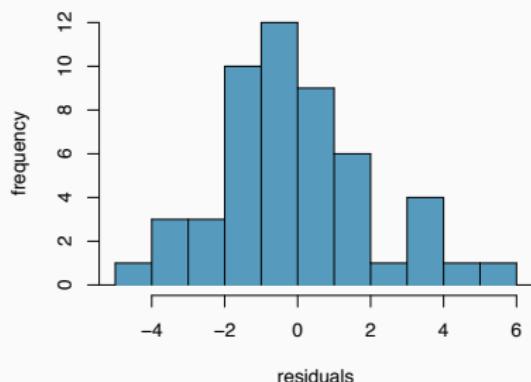
$$\% \text{ HS grad} = 86 \quad \% \text{ in poverty} = 16.8$$

$$\widehat{\% \text{ in poverty}} = 64.68 - 0.62 * 86 = 11.36$$

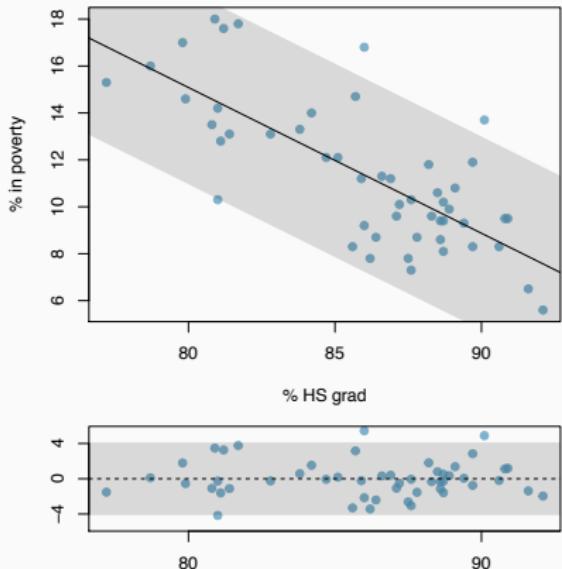
$$\begin{aligned} e &= \% \text{ in poverty} - \widehat{\% \text{ in poverty}} \\ &= 16.8 - 11.36 = 5.44 \end{aligned}$$

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

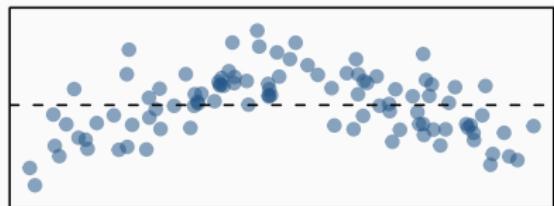
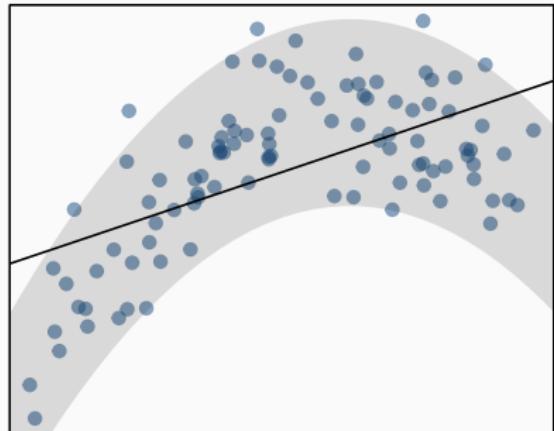


- * The variability of points around the least squares line should be roughly constant.
- * This implies that the variability of residuals around the 0 line should be roughly constant as well.
- * Also called [homoscedasticity](#).
- * Check using a histogram or normal probability plot of residuals.

CHECKING CONDITIONS

What condition is this linear model obviously violating?

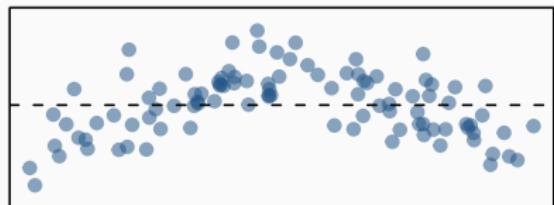
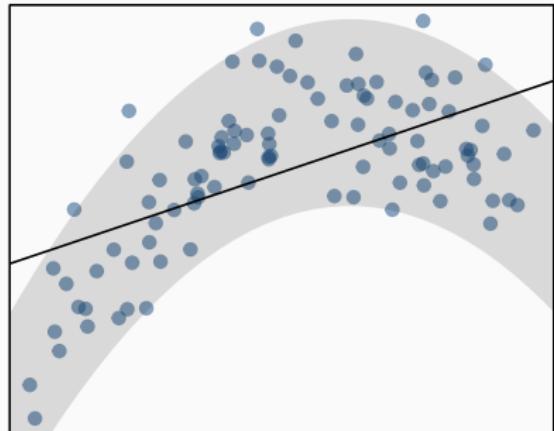
- * Constant variability
- * Linear relationship
- * Normal residuals
- * No extreme outliers



CHECKING CONDITIONS

What condition is this linear model obviously violating?

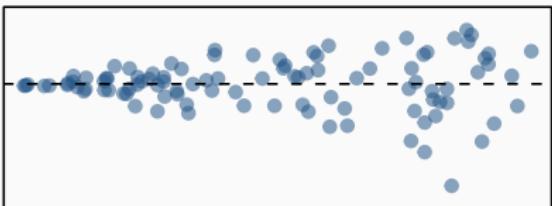
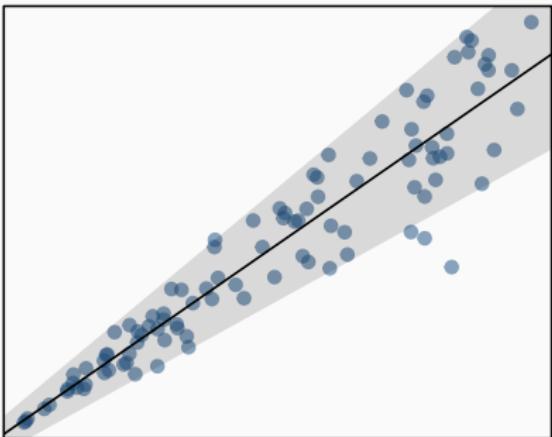
- * Constant variability
- * **Linear relationship**
- * Normal residuals
- * No extreme outliers



CHECKING CONDITIONS

What condition is this linear model obviously violating?

- * Constant variability
- * Linear relationship
- * Normal residuals
- * No extreme outliers

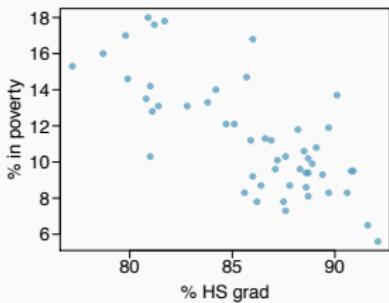


- * The strength of the fit of a linear model is most commonly evaluated using R^2 .
- * R^2 is calculated as the square of the correlation coefficient.
- * It tells us what percent of variability in the response variable is explained by the model.
- * The remainder of the variability is explained by 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 R²

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

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

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

- * 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

- * northeast
- * midwest
- * west
- * south
- * 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

- * northeast
- * midwest
- * west
- * south
- * cannot tell

Line fitting, residuals, and correlation

Fitting a line by least squares regression

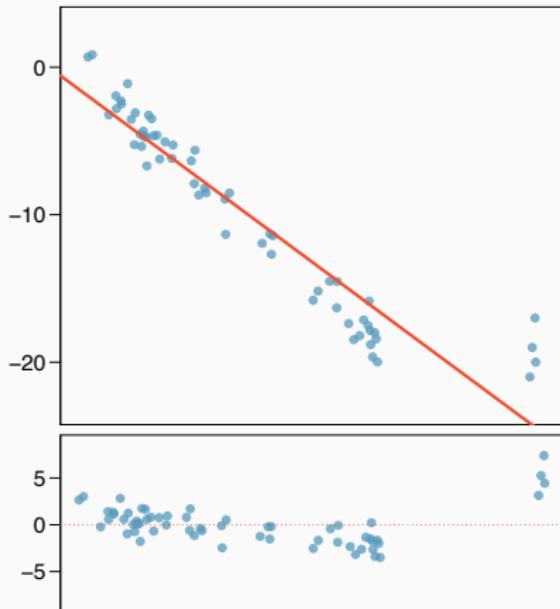
Types of outliers in linear regression

Inference for linear regression

TYPES OF OUTLIERS

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

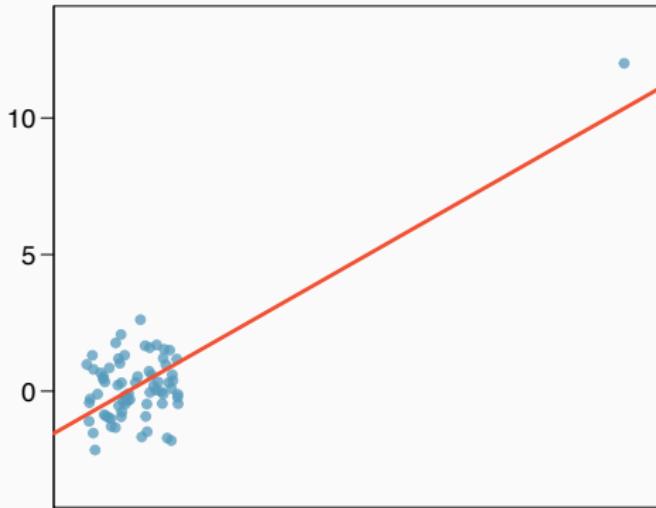
To answer this question think of where the regression line would be with and without the outlier(s). Without the outliers the regression line would be steeper, and lie closer to the larger group of observations. With the outliers the line is pulled up and away from some of the observations in the larger group.



TYPES OF OUTLIERS

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

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

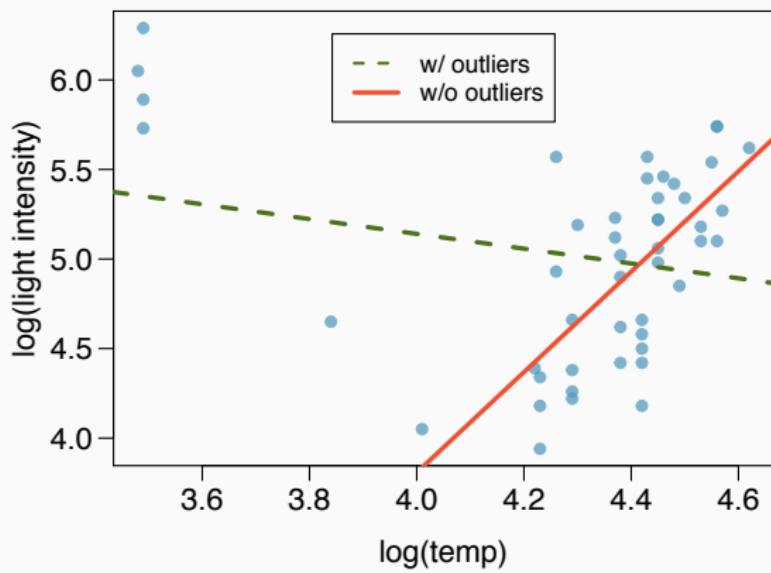


SOME TERMINOLOGY

- * Outliers are points that lie away from the cloud of points.
- * 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 itnot an influential point.

INFLUENTIAL POINTS

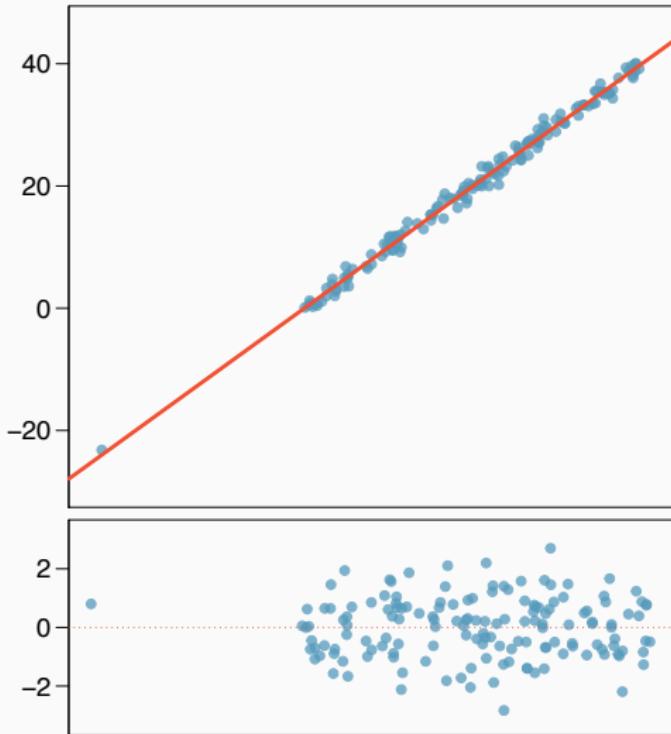
Data are available on the log of the surface temperature and the log of the light intensity of 47 stars in the star cluster CYG OB1.



TYPES OF OUTLIERS

Which of the below best describes the outlier?

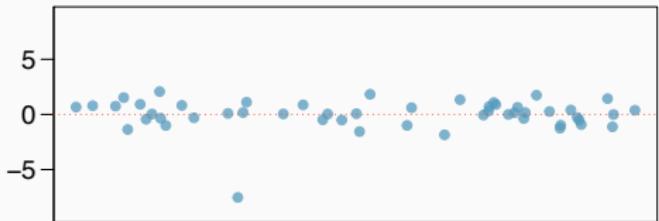
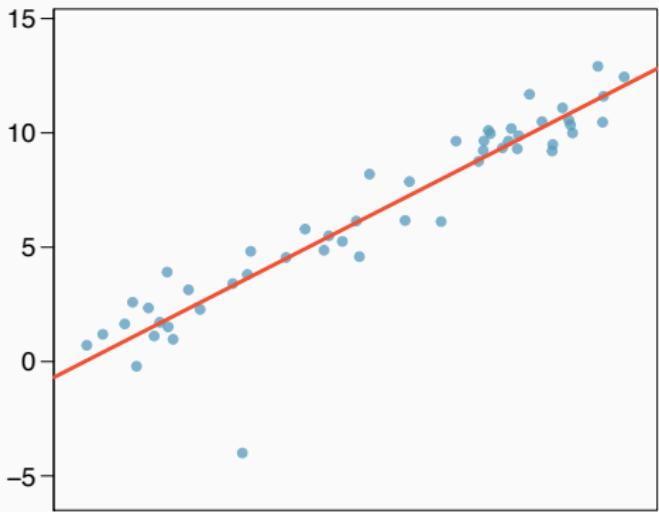
- * influential
- * high leverage
- * none of the above
- * there are no outliers



TYPES OF OUTLIERS

Does this outlier influence the slope of the regression line?

Not much...

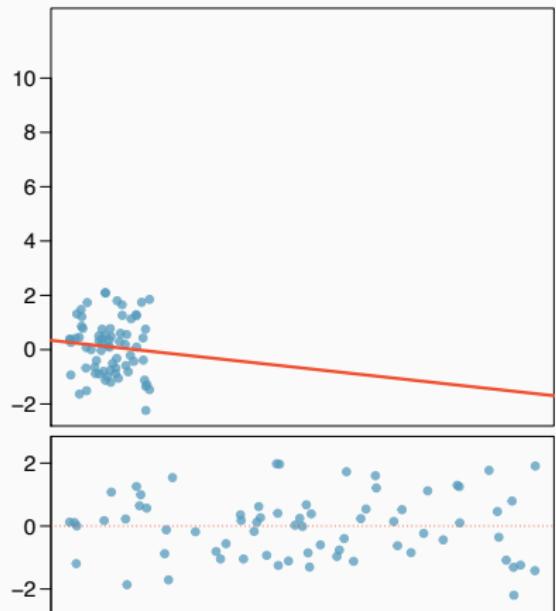


Which of following is true?

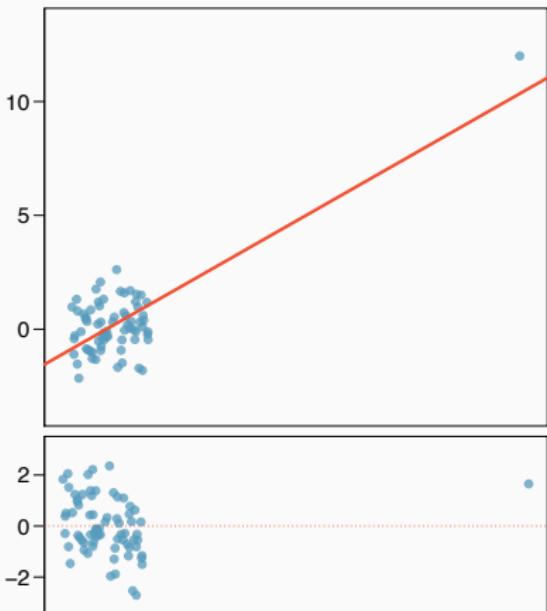
- * Influential points always change the intercept of the regression line.
- * Influential points always reduce R^2 .
- * It is much more likely for a low leverage point to be influential, than a high leverage point.
- * When the data set includes an influential point, the relationship between the explanatory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

RECAP (CONT.)

$R = 0.08, R^2 = 0.0064$



$R = 0.79, R^2 = 0.6241$



Line fitting, residuals, and correlation

Fitting a line by least squares regression

Types of outliers in linear regression

Inference for linear regression

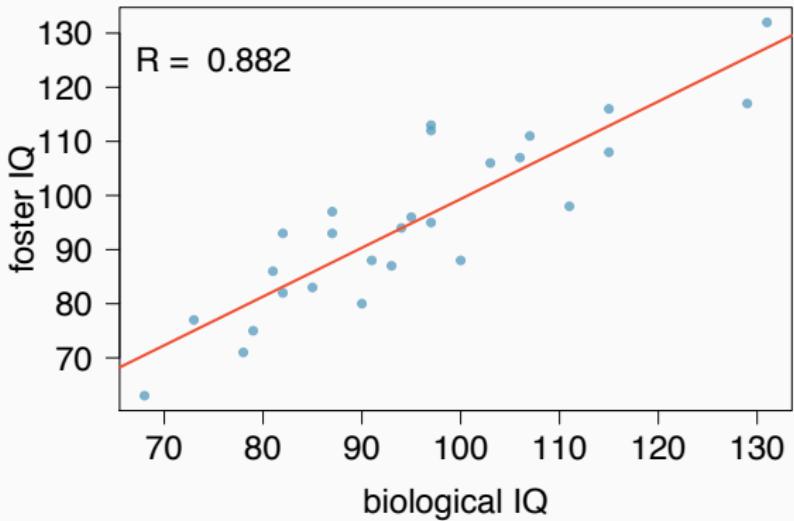
Understanding regression output from software

HT for the slope

CI for the slope

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.



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

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

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?

- * $H_0 : b_0 = 0; H_A : b_0 \neq 0$
- * $H_0 : \beta_0 = 0; H_A : \beta_0 \neq 0$
- * $H_0 : b_1 = 0; H_A : b_1 \neq 0$
- * $H_0 : \beta_1 = 0; H_A : \beta_1 \neq 0$

TESTING FOR THE SLOPE (CONT.)

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

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

Remember: Test statistic, $T = \frac{\text{point estimate} - \text{null value}}{\text{SE}}$

- * 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

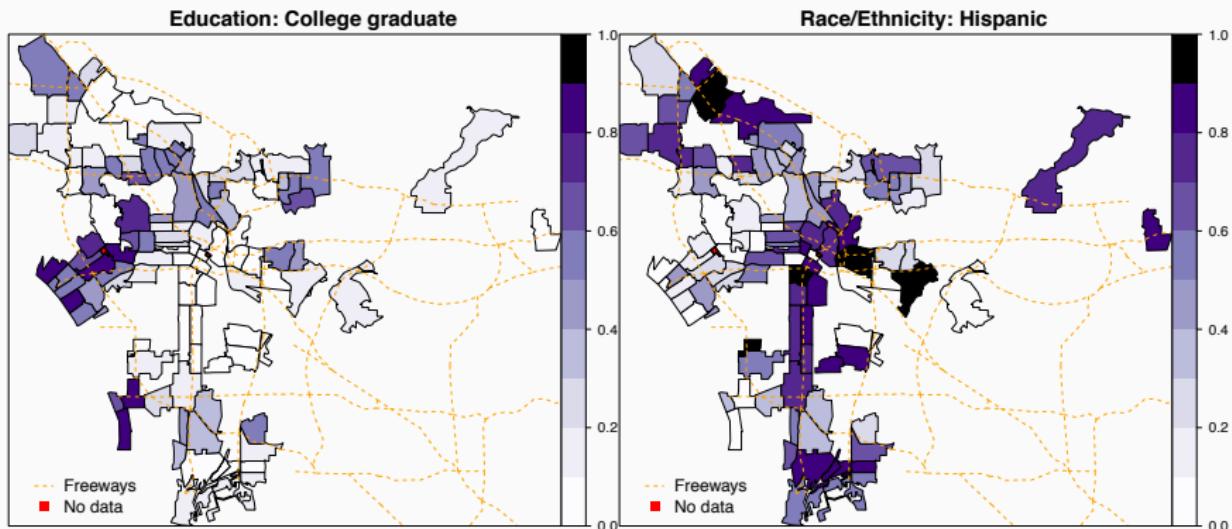
$$T = \frac{0.9014 - 0}{0.0963} = 9.36$$

$$df = 27 - 2 = 25$$

$$p\text{-value} = P(|T| > 9.36) < 0.01$$

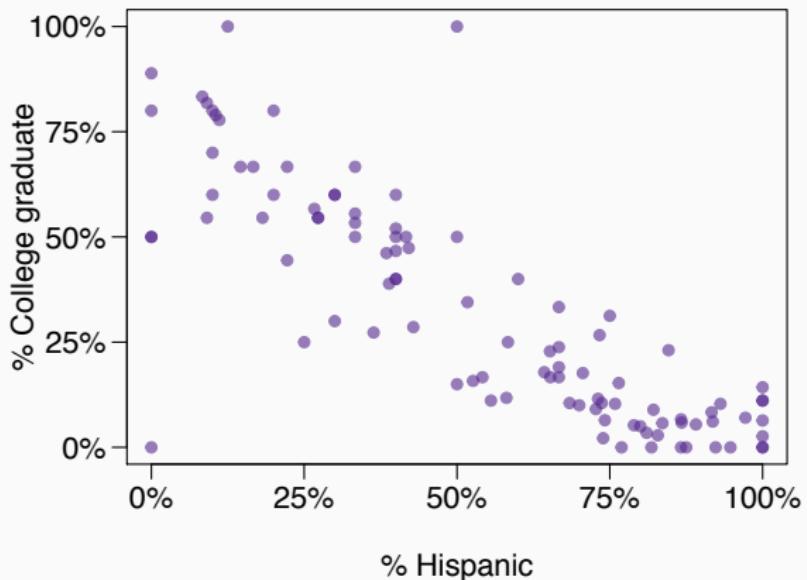
% COLLEGE GRADUATE VS. % HISPANIC IN LA

What can you say about the relationship between % college graduate and % Hispanic in a sample of 100 zip code areas in LA?



% COLLEGE EDUCATED VS. % HISPANIC IN LA - ANOTHER LOOK

What can you say about the relationship between % college graduate and % Hispanic in a sample of 100 zip code areas in LA?



% COLLEGE EDUCATED VS. % HISPANIC IN LA - LINEAR MODEL

Which of the below is the best interpretation of the slope?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.7290	0.0308	23.68	0.0000
%Hispanic	-0.7527	0.0501	-15.01	0.0000

- * A 1% increase in Hispanic residents in a zip code area in LA is associated with a 75% decrease in % of college grads.
- * A 1% increase in Hispanic residents in a zip code area in LA is associated with a 0.75% decrease in % of college grads.
- * An additional 1% of Hispanic residents decreases the % of college graduates in a zip code area in LA by 0.75%.
- * In zip code areas with no Hispanic residents, % of college graduates is expected to be 75%.

% COLLEGE EDUCATED VS. % HISPANIC IN LA - LINEAR MODEL

Do these data provide convincing evidence that there is a statistically significant relationship between % Hispanic and % college graduates in zip code areas in LA?

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.7290	0.0308	23.68	0.0000
hispanic	-0.7527	0.0501	-15.01	0.0000

Yes, the p-value for % Hispanic is low, indicating that the data provide convincing evidence that the slope parameter is different than 0.

How reliable is this p-value if these zip code areas are not randomly selected?

Not very...

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

- * $9.2076 \pm 1.65 \times 9.2999$ $n = 27$ $df = 27 - 2 = 25$
- * $0.9014 \pm 2.06 \times 0.0963$ $95\% : t_{25}^* = 2.06$
- * $0.9014 \pm 1.96 \times 0.0963$ $0.9014 \pm 2.06 \times 0.0963$
- * $9.2076 \pm 1.96 \times 0.0963$ $(0.7, 1.1)$

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

- * Hypothesis test:

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

- * Confidence interval:

$$b_1 \pm t_{df=n-2}^* \text{SE}_{b_1}$$

- * 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.
- * 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.

Hw 4:

Simple Linear Model (population)

Data	Population
$(X_i, Y_i) \xrightarrow{i.i.d.} X_i \in \mathbb{R}^d$	$(X, Y) \sim P_0$

unknown distribution of the pair

Y is capitalized to denote that it is a random variable. x is fixed (not random)

Simple Linear Model : $\underbrace{\mathbb{E}[Y|x]}_{\text{conditional Expectation}} = \underbrace{a + bx}_{\substack{\text{linear in } x \\ \text{placeholder for all possible } x \text{ in } X}} = \mathbb{E}[Y|X=x]$

Gaussian Measurement Model (data)

$$Y_i = a + bX_i + \epsilon_i \quad \leftarrow \text{bridge between data \& population}$$

$$\epsilon_i \xrightarrow{i.i.d.} N(0, \sigma^2) \quad [\sigma^2 \text{ unknown}]$$

$\epsilon_i \leftarrow$ we can only estimate, not observable \leftarrow model Error

$$(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$$

Least Square $(\hat{a}, \hat{b}) = \underset{a, b}{\text{argmin}} \sum_{i=1}^n (Y_i - a - bX_i)^2$

$$\hat{a} = \frac{(\sum_{i=1}^n X_i) \cdot (\sum_{i=1}^n Y_i) - (\sum_{i=1}^n X_i)(\sum_{i=1}^n X_i Y_i)}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2}$$

$$\hat{b} = \frac{n \cdot \sum_{i=1}^n (X_i Y_i) - (\sum_{i=1}^n X_i)(\sum_{i=1}^n Y_i)}{n \cdot \sum_{i=1}^n (X_i)^2 - (\sum_{i=1}^n X_i)^2}$$

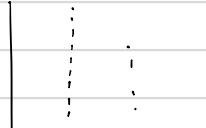
weighted sums of recorded responses

For some weights

$$\begin{aligned} \hat{a} &= \sum_{i=1}^n w_i Y_i \\ \hat{b} &= \sum_{i=1}^n w_i X_i \end{aligned} \quad \left. \begin{array}{l} \Rightarrow \text{linear functions of } Y_i \\ \Rightarrow \text{linear functions of } \epsilon_i \\ \text{(although we dont observe } \epsilon_i \text{)} \end{array} \right.$$

Example 1 :

can never have
(not allowed)
 \Rightarrow violate the models



Example 2 :

$$\hat{b} = 0$$

shows that Y is uncorrelated with X



LS is unbiased

- Under Gaussian Measurement Model (GMM) $E[\hat{a}] = a$
(GMM, $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$) $E[\hat{b}] = b$

- Residuals: $\hat{\epsilon}_i = y_i - (\hat{a} - \hat{b}x_i)$ \rightarrow fitted value
 $\Rightarrow E[\hat{\epsilon}_i] = \epsilon_i$
 $\hat{\epsilon}_i$ are unbiased as well

$\text{Var}(\hat{\epsilon}_i) \neq \sigma^2$, $\text{Var}(\hat{\epsilon}_i)$ depends on σ^2 (even if the model is perfect)
 $\sigma^2 = ?$

$\text{Var}(\text{standardized residual}_i) = \sigma^2$
 $\hookrightarrow E[\text{standardized residual}] = \epsilon_i$ unbiased
We look at standardized residuals to do model checking.

- $\hat{\sigma}^2 = ?$
 $\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\epsilon}_i^2$
of regression parameters (a and b)

What if # of parameter = n?

$$E[Y | \vec{x}] = a + b_1 x_1 + b_2 x_2 + \dots + b_{n-1} x_{n-1}$$

$\hat{\sigma}^2 \rightarrow \infty$ (have large predicting error) $\leftarrow Y_i - \hat{Y}_i$ is huge \Rightarrow overfitting

$$p = O(\sqrt{n}) \quad \frac{p}{\sqrt{n}} \xrightarrow{n \rightarrow \infty} 0 \quad \begin{matrix} \text{(regression)} \\ p \Rightarrow \# \text{ of parameters in your model} \end{matrix}$$

$$\Rightarrow \dim(x) + 1$$

Model Misfit:

An alternative to a linear model is a polynomial model

Quadratic Model $E[Y | x] = c + dx + ex^2$ (c and d and e are unknown)

\Leftrightarrow Two-variable Linear Model

$$\text{Call } u = x^2$$

$$E[Y | X, U] = c + dx + eu$$

Misfit

fit a linear model on a quadratic case

$$(*) Y_i = c + dx_i + eu_i + \epsilon_i$$

If using \hat{a}, \hat{b} from least square as before,
 $\Rightarrow \hat{b}$ is no longer unbiased $\Rightarrow \hat{b}$ is biased

$$E[\hat{b}] = d + e \underbrace{\frac{n \cdot \sum_{i=1}^n (X_i U_i) - (\sum_{i=1}^n X_i) (\sum_{i=1}^n U_i)}{n \sum_{i=1}^n (X_i^2) - (\sum_{i=1}^n X_i)^2}}$$

(X^2)

only equals 0 when X_i is uncorrelated with $U_i \Rightarrow$ must be correlated

\Rightarrow cannot be 0

If two variables are not correlated

- \rightarrow just include one variable in the estimate (delete one)
- \Rightarrow unbiased estimator \hat{b} ($Y_i = c + d X_i + \epsilon_i$)

5/17/2018

Replicate Measurements

→ g distinct values of explanatory variable x

→ For each x , 10 replicate measurements

$$Y_{ij} = \alpha + \beta x_i + \epsilon_{ij}$$

$$i = 1, \dots, m$$

m = # of distinct values

$$j = 1, \dots, k$$

Y_{ij} = j^{th} measurement taken at X_i

We will assume ϵ_{ij} is uncorrelated for all i and j

→ Use replicate measurement to estimate σ^2 , that does not rely too much on the model

→ If the model is fitted incorrectly, the residuals include measurement error from ϵ_{ij} as well as model misfit.

Suppose $Y_{ij} = c + d X_i + e u_i + \epsilon_{ij}$, but simple linear model is fit (by least squares)

Then, $Y_{11}, Y_{12}, \dots, Y_{1k}$ } replicates are k uncorrelated random variables
r.v.

$$E[Y_{11}] = E[c + d X_1 + e u_1 + \epsilon_{11}] = c + d X_1 + e u_1$$

$$\text{Var}(Y_{11}) = \sigma^2$$

$$S^2 = \frac{1}{k-1} \sum_{j=1}^k (Y_{1j} - \bar{Y}_1)^2$$

$$\bar{Y}_1 = \frac{1}{k} \sum_{j=1}^k Y_{1j}$$

$$S_1^2, S_2^2, \dots, S_m^2 \text{ (no regression)}$$

$$\text{Final } S^2_{\text{pooled}} = \frac{1}{m} \sum_{i=1}^m S_i^2 : \text{def } S_p^2$$

$$S_m^2 = \frac{1}{k-1} \sum_{j=1}^k (Y_{kj} - \bar{Y}_k)^2$$

← prediction interval is shorter

Is the model misfit? (the statistical way)

Residual sum of squares : (RSS)

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^k (Y_{ij} - \hat{Y}_{ij})^2 \leftarrow \hat{Y}_i \text{ not } \hat{Y}_{ij} \text{ because we are only predicting } X_i \\ & = \underbrace{\sum_{i=1}^n \sum_{j=1}^k (Y_{ij} - \bar{Y}_i)^2}_{(\text{MESS}) \text{ measurement error}} + K \cdot \underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y}_i)^2}_{\text{Model Misfit (MMSS)}} \end{aligned}$$

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i$$

MESS has $m(k-1)$ degrees of freedom

$$\text{MESS} \sim \chi_{df}^2$$

$$\text{MMSS} \sim \chi_{m-2}^2$$

If $E_{ij} \sim N(0, \sigma^2)$ and there is no model misfit
 Design a test statistic to check if there is a model misfit.

$$\frac{k \cdot \text{MSESS} / (m-2)}{\text{MESS} / (m \cdot (k-1))} \sim F_{m-2, m(k-1)}$$

(Fisher's distribution)

↑ If null hypothesis is true, use this to reject

When to reject?

↳ We reject = there is a model misfit = large values of $F (> 3, > 5)$

Confidence and prediction bounds

Predicted Value : $\hat{Y} = \hat{a} + \hat{b}x$

time value of function for Y



$$\text{Var}(\hat{Y} - Y) = \text{Var}(\hat{a} + \hat{b}x - Y) = \text{Var}(\hat{a} + \hat{b}x - a - bx - \varepsilon) = \text{Var}(\hat{a} - a) + (\hat{b} - b)x - \varepsilon$$

$$\text{Var}(\hat{Y} - Y) = \sigma^2 \left(1 + \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})^2\right)$$

Prediction intervals (L, R)

$$L = (\hat{a} + \hat{b}x) - z_{1-\alpha/2} \hat{\sigma} \sqrt{1 + \frac{1}{m} + \frac{(x - \bar{x})^2}{\sum_{i=1}^m (x_i - \bar{x})^2}}$$

↑ quantile of standard normal

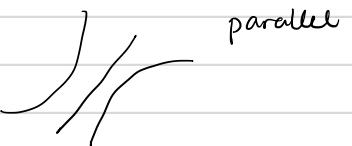
⇒ for measurement error

$$R = (\hat{a} + \hat{b}x) + z_{1-\alpha/2} \hat{\sigma} \sqrt{1 + \frac{1}{m} + \frac{(x - \bar{x})^2}{\sum_{i=1}^m (x_i - \bar{x})^2}}$$

↑ quantile of standard normal

Supposed to be curved, not parallel

$$m \text{ small } z_{1-\alpha/2} \rightarrow t_{df, 1-\alpha/2} \text{ student } df = m-2$$



* contrast ggplot2's parameter with the pooled estimator

If x is far from \bar{x} , then PI is wider.



CI : $(L, R) \leftarrow \text{for } E[Y|x]$ prediction interval for Y

$$L = (\hat{a} + \hat{b}x) - t_{m-2, 1-\alpha/2} \cdot \sqrt{\hat{\sigma}^2 \frac{1}{m} + \frac{(x - \bar{x})^2}{\sum_{i=1}^m (x_i - \bar{x})^2}}$$

$$R = (\hat{a} + \hat{b}x) + t_{m-2, 1-\alpha/2} \cdot \sqrt{\hat{\sigma}^2 \frac{1}{m} + \frac{(x - \bar{x})^2}{\sum_{i=1}^m (x_i - \bar{x})^2}}$$

$H_0: b=0$ (significance test)

(Is x linearly associated to Y)

$$T = \frac{\hat{b}}{\text{SE}(\hat{b})} \quad \text{SE}(\hat{b}) = \text{standard error} \sim \text{student } t_{n-2}$$

Reject H_0 , $|T| > T_{\text{observed}}$ p-value = $2P(T > T_{\text{observed}})$

reject p-value < α