

# Statistical Inference

## Introduction to Linear Regression

*Behnam Bahrak*  
*Spring 2020*

1 of 23



## Introduction to Linear Regression

- So far we worked with:
  - single numerical variable
  - single categorical variable
  - relationship between a numerical and a categorical variable
  - relationship between two categorical variables
- Linear Regression:
  - the relationship between two numerical variables.
  - Correlation: a measure of the strength of the linear relationship between two numerical variables



Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

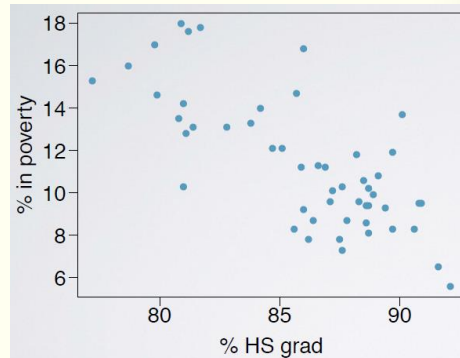


2 of 23



## Poverty vs. HS Grad Rate

- Data: 50 states + DC
- Poverty line in US: income below \$23,050 for a family of 4 in 2012
- Response?  $y$   
% in poverty
- Explanatory?  $x$   
% HS grad
- Relationship?  
linear, negative,  
moderately strong



Statistical Inference

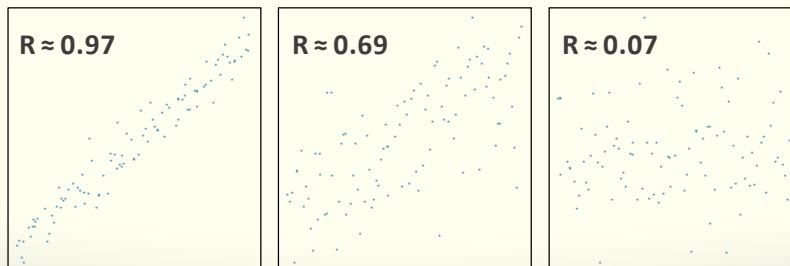
Behnam Bahrak  
bahrak@ut.ac.ir

&lt; 3 of 23 &gt;

## Correlation

- Describes the strength of the **linear association** between two variables and is denoted as  $R$
- **Property 1.** The magnitude (absolute value) of the correlation coefficient measures the strength of the **linear association** between two numerical variables

$$R = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$



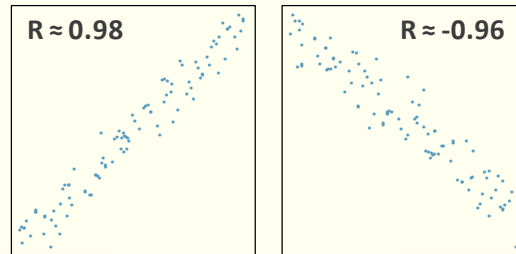
Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

&lt; 4 of 23 &gt;

## Properties of Correlation

- **Property 2.** The sign of the correlation coefficient indicates the direction of association



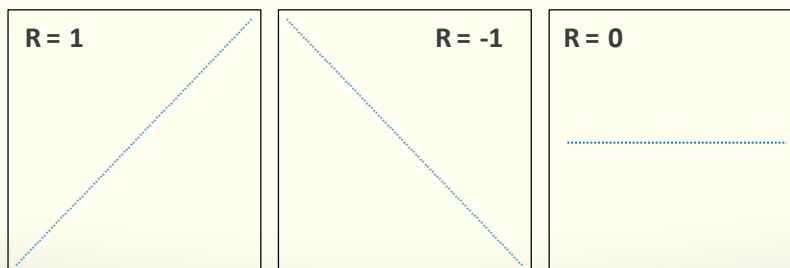
Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

5 of 23

## Properties of Correlation

- **Property 3.** The correlation coefficient is always between -1 (perfect negative linear association) and 1 (perfect positive linear association)
- $R = 0$  indicates **no linear relationship**



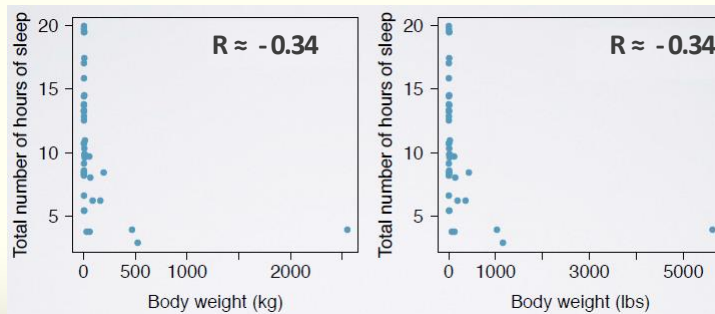
Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

6 of 23

## Properties of Correlation

- **Property 4.** The correlation coefficient is unitless, and is not affected by changes in the center or scale of either variable (such as unit conversions)



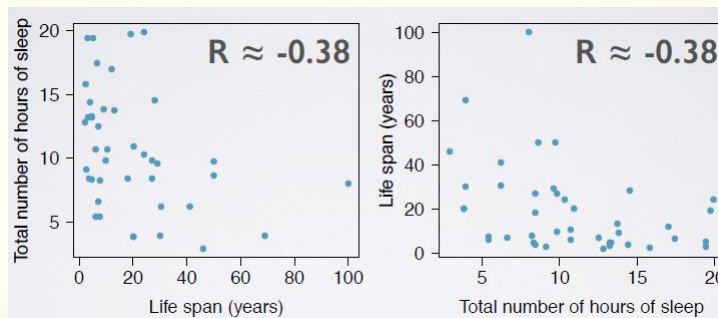
Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

&lt; 7 of 23 &gt;

## Properties of Correlation

- **Property 5.** The correlation of  $X$  with  $Y$  is the same as of  $Y$  with  $X$ .



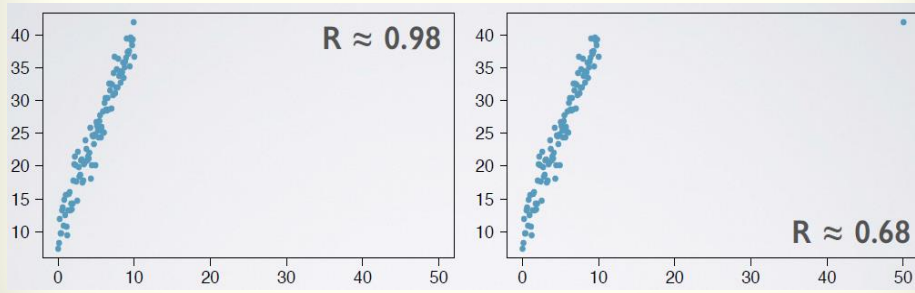
Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

&lt; 8 of 23 &gt;

## Properties of Correlation

- **Property 6.** The correlation coefficient is sensitive to outliers.



Statistical Inference

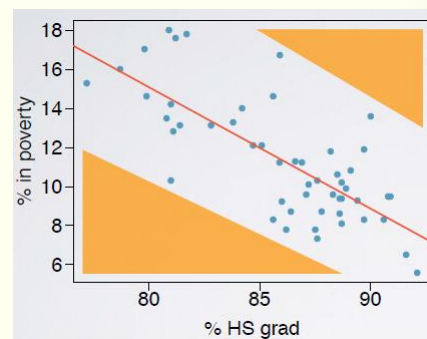
Behnam Bahrak  
bahrak@ut.ac.ir

9 of 23

## Question

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

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



Statistical Inference

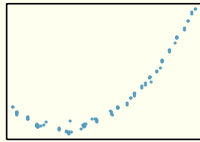
Behnam Bahrak  
bahrak@ut.ac.ir

10 of 23

## Question

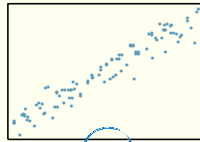
- Which of the following has the strongest correlation, i.e. correlation coefficient closest to +1 or -1?

very strong,  
but not linear



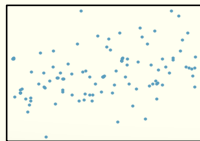
(a)

strongest  
linear



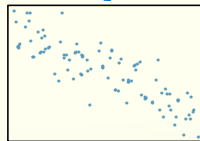
(b)

even weaker



(c)

weaker



(d)



Statistical Inference

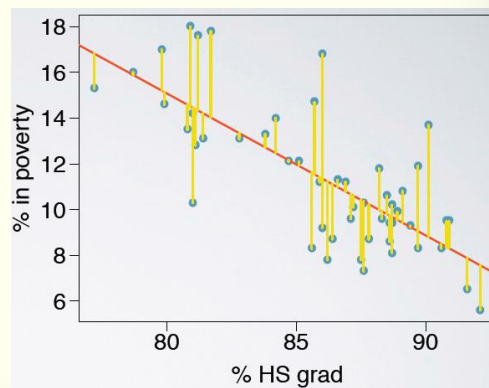
Behnam Bahrak  
bahrak@ut.ac.ir

11 of 23

## Residuals

- Leftovers from the model fit
- **Data = Fit + Residual**
- Difference between the observed and predicted  $y$

$$\text{residual: } e_i = y_i - \hat{y}_i$$

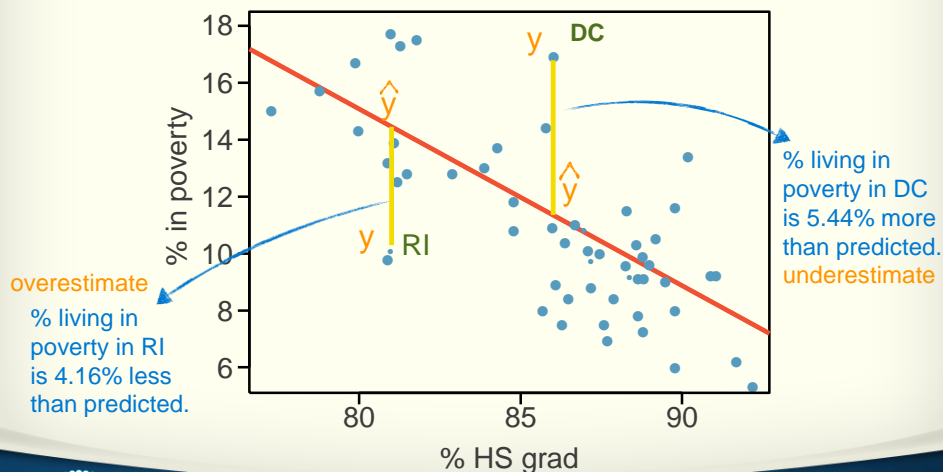


Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

12 of 23

## Residuals



Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

13 of 23

## A measure for the best line

**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

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



Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

14 of 23

## 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 more than twice as bad

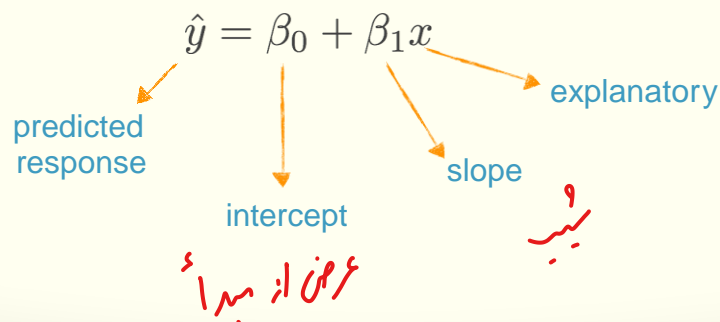


Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

15 of 23

## Least Square Line



Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

16 of 23



## Notation

معادله رگرسیونی

	parameter	point estimate
intercept	$\beta_0$	$b_0$
slope	$\beta_1$	$b_1$



Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

17 of 23

## Estimating the regression parameters: slope

slope:

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

$s_x$  : SD of  $x$

$s_y$  : SD of  $y$

$R = \text{cor}(x, y)$



Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

18 of 23

## Example

- The standard deviation of % living in poverty is 3.1%, and the standard deviation of % HS graduates is 3.73%. Given that the correlation between these variable is -0.75, what is the slope of the regression line for predicting % living poverty from % HS graduates?

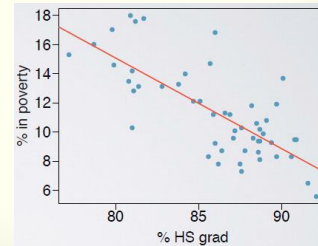
$$S_y = 3.1\%$$

$$S_x = 3.73\%$$

$$R = -0.75$$

$$b_1 = \frac{S_y}{S_x} \times R = \frac{3.1}{3.73} \times (-0.75) = -0.62$$

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



Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

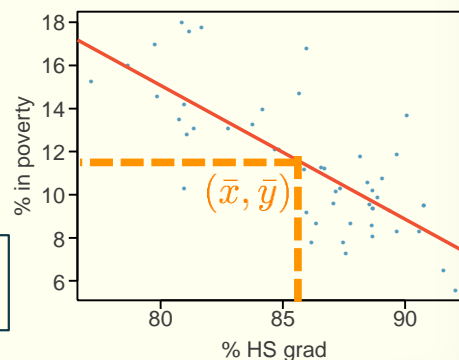
19 of 23

## Estimating the regression parameters: intercept

- The least squares line always goes through  $(\bar{x}, \bar{y})$

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

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



Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir

20 of 23

## Example

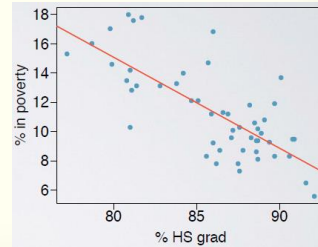
- Given that the average % living in poverty is 11.35%, and the average % HS graduates is 86.01%, what is the intercept of the regression line for predicting % living poverty from % HS graduates?

$$\bar{y} = 11.35\%$$

$$\bar{x} = 86.01\%$$

$$b_0 = \bar{y} - b_1\bar{x} = 11.35 - (-0.62)(86.01) = 64.68$$

States with no HS graduates are expected on average to have 64.68% of their residents living below the poverty line.



Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir



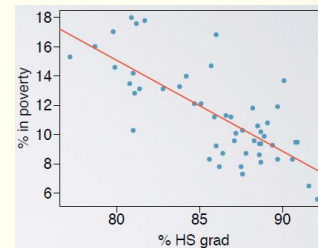
21 of 23



## Example

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

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	64.78	6.80	9.52	0.00
hsgrad	-0.62	0.08	-7.86	0.00



Statistical Inference

Behnam Bahrak  
bahrak@ut.ac.ir



22 of 23



## Recap

- **intercept:** When  $x = 0$ ,  $y$  is expected to equal the intercept.
  - may be meaningless in context of the data, and only serve to adjust the height of the line
- **slope:** For each unit increase in  $x$ ,  $y$  is expected to be higher/lower on average by the slope.

