

---

# THE LINEAR REGRESSION PROCESS

---

FROM GENERATION TO INTERPRETATION

PRESENTED BY

FRANK BRIODY  
*frankbriody@gmail.com*

*Prospect High School  
Mt. Prospect, IL*



2022 MMC CONFERENCE OF WORKSHOPS



*Photo credit - Vito Palmisano*

## Contents

---

<b>1</b>	<b>Standard Deviation: The Non-Resistant Measure of Spread</b>	<b>2</b>
<b>2</b>	<b>Getting Least Squares Line of Best Fit</b>	<b>3</b>
2.1	From Summary Statistics . . . . .	3
2.2	From Output . . . . .	4
2.3	Interpretation . . . . .	4
2.3.1	Slope . . . . .	4
2.3.2	Y-Intercept . . . . .	4
2.4	Predicted Values and Residuals . . . . .	5
2.5	The Coefficient of Determination $r^2$ - Comparing Models . . . . .	6
<b>3</b>	<b>Summary &amp; Examples</b>	<b>7</b>
3.1	Getting $r$ . . . . .	7
3.2	Five R Properties . . . . .	7

## The Story

A statistics teacher gives a quiz to a class. The scores were 2, 4, 6, 8, and 15 with one student being absent. Absent student returns the next day...

**Student:** How am I going to do on the quiz?

**Teacher:** Well, the class average was...

## 1 Standard Deviation: The Non-Resistant Measure of Spread

How much variability, on average, is there around the mean?

<div>L1</div> <div>1</div> <div>2</div> <div>3</div> <div>4</div> <div>5</div> <div>-----</div> <div>L3(5) = 1</div>	<div>L2</div> <div>2</div> <div>4</div> <div>6</div> <div>8</div> <div>15</div> <div>-----</div>	<div>L3</div> <div>0</div> <div>0</div> <div>0</div> <div>1</div> <div>-----</div>	<div>EDIT CALC TESTS</div> <div>1:1-Var Stats</div> <div>2:2-Var Stats</div> <div>3:Med-Med</div> <div>4:LinReg(ax+b)</div> <div>5:QuadReg</div> <div>6:CubicReg</div> <div>7:QuartReg</div>	<div>1-Var Stats</div> <div>List:L2</div> <div>FreqList:</div> <div>Calculate</div>	<div>1-Var Stats</div> <div><math>\bar{x}=7</math></div> <div><math>\Sigma x=35</math></div> <div><math>\Sigma x^2=345</math></div> <div><math>Sx=5</math></div> <div><math>\sigma x=4.472135955</math></div> <div><math>n=5</math></div>	<div>1-Var Stats</div> <div><math>n=5</math></div> <div><math>\min X=2</math></div> <div><math>Q1=3</math></div> <div><math>\text{Med}=6</math></div> <div><math>Q3=11.5</math></div> <div><math>\max X=15</math></div>
--	--	--	--	---	---	---

Score    Deviation    Squared Deviation

$x$

2

4

6

8

15

## The Story Part 2

A statistics teacher gives a quiz to a class. The scores were 2, 4, 6, 8, and 15 with one student being absent. After surveying the class, the teacher knows the hours studied were 1, 2, 3, 4 and 5, respectively. Absent student returns the next day...

**Student:** How am I going to do on the quiz?

**Teacher:** That depends - how long did you study?

## 2 Getting Least Squares Line of Best Fit

---

Hours ( $x$ )	1	2	3	4	5
Score ( $y$ )	2	4	6	8	15

### 2.1 From Summary Statistics

Formulas (given):  $\hat{y} = a + bx$        $b = r \frac{s_y}{s_x}$        $a = \bar{y} - b\bar{x}$

Descriptive Statistics: x, y

Variable	N	N*	Mean	SE Mean	StDev	Minimum	Q1	Median	Q3	Maximum
x	5	0	3.000	0.707	1.581	1.000	1.500	3.000	4.500	5.000
y	5	0	7.00	2.24	5.00	2.00	3.00	6.00	11.50	15.00

Correlations: x, y

Pearson correlation of x and y = 0.949

P-Value = 0.014

## 2.2 From Output

### Regression Analysis: y versus x

The regression equation is

$$y = -2.00 + 3.00 x$$

Predictor	Coef	SE Coef	T	P
Constant	-2.000	1.915	-1.04	0.373
x	3.0000	0.5774	5.20	0.014

s = 1.82574      R-sq = 90.0%      R-Sq(adj) = 86.7%

### Analysis of Variance

Source	DF	SS	MS	F	P
Regression	1	90.000	90.000	27.00	0.014
Residual Error	3	10.000	3.333		
Total	4	100.000			

## 2.3 Interpretation

### 2.3.1 Slope

Slope represents the **predicted** change in response associated with each unit increase in the explanatory variable, **on average**.

### 2.3.2 Y-Intercept

Y-intercept is the predicted value when the explanatory ( $x$ ) is 0. [Often the y-intercept is useless.]

## 2.4 Predicted Values and Residuals

$$\hat{y} = -2 + 3x$$

Hours	Score	Predicted	Residual
$x$	$y$	$\hat{y}$	$y - \hat{y}$
1	2		
2	4		
3	6		
4	8		
5	15		

- Predicted  $\hat{y}$ : substitute explanatory ( $x$ ) values into regression equation.
- Residual  $y - \hat{y}$  (also called *regression error*); actual minus predicted.

### The Story Part 3

A statistics teacher gives a quiz to a class. The scores were 2, 4, 6, 8, and 15 with one student being absent. After surveying the class, the teacher knows the hours studied were 1, 2, 3, 4 and 5, respectively. Absent student returns the next day...

**Student:** How am I going to do on the quiz?

**Teacher:** That depends - how long did you study?

**Student:** Does how long I studied really make a difference?

## 2.5 The Coefficient of Determination $r^2$ - Comparing Models

**Regression Analysis: y versus x**

The regression equation is

$$y = -2.00 + 3.00x$$

Predictor	Coef	SE Coef	T	P
Constant	-2.000	1.915	-1.04	0.373
x	3.0000	0.5774	5.20	0.014

s = 1.82574      R-sq = 90.0%      R-Sq(adj) = 86.7%

**Analysis of Variance**

Source	DF	SS	MS	F	P
Regression	1	90.000	90.000	27.00	0.014
Residual Error	3	10.000	3.333		
Total	4	100.000			

Hours	Score	Predicted	Error	(Error) <sup>2</sup>	Residual	(Residual) <sup>2</sup>
$x$	$y$	$\hat{y}$	$y - \bar{y}$	$(y - \bar{y})^2$	$y - \hat{y}$	$(y - \hat{y})^2$
1	2					
2	4					
3	6					
4	8					
5	15					

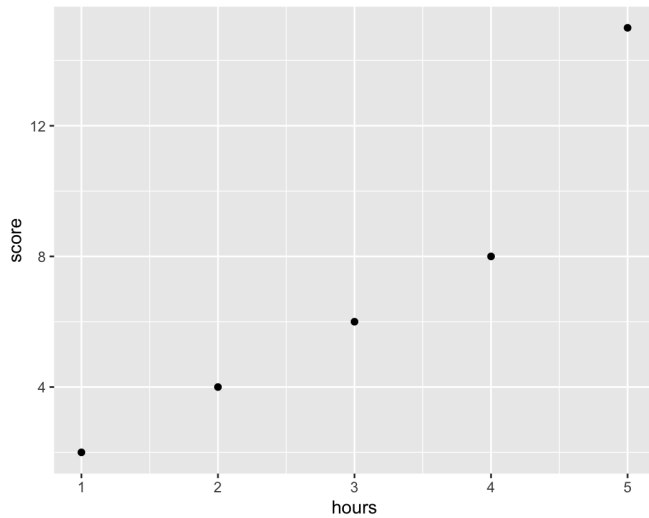
### 3 Summary & Examples

---

#### 3.1 Getting $r$

- $r = \frac{\sum z_x \cdot z_y}{n-1}$
- Never calculate by hand; use calculator or computer output.
- Know formula *properties*.

The r Formua



```
Link9
y=a*x+b
a=3
b=-2
r^2=.9
r=.9486832981
```

Correlations: x, y  
Pearson correlation of x and y = 0.949  
P-Value = 0.014

#### 3.2 Five R Properties

- Examples