

Unit 12 Simple Linear Regression and Correlation

*“Assume that a statistical model such as a linear model
is a good first start only”*

- Gerald van Belle

Is higher blood pressure in the mom associated with a lower birth weight of her baby? Simple linear regression explores the relationship of **one continuous outcome** (Y=birth weight) with **one continuous predictor** (X=blood pressure). At the heart of statistics is the fitting of models to data followed by an examination of how the models perform.

-1- **“somewhat useful”**

A fitted model is somewhat useful if it permits exploration of hypotheses such as “higher blood pressure during pregnancy is associated with statistically significant lower birth weight” and it permits assessment of confounding, effect modification, and mediation. These are ideas that will be developed further in BIOSTATS 640 Unit 5, ***Normal Theory Regression***.

-2- **“more useful”**

The fitted model is more useful if it can be used to predict the outcomes of future observations. For example, we might be interested in predicting the birth weight of the baby born to a mom whose systolic blood pressure is 145 mm Hg.

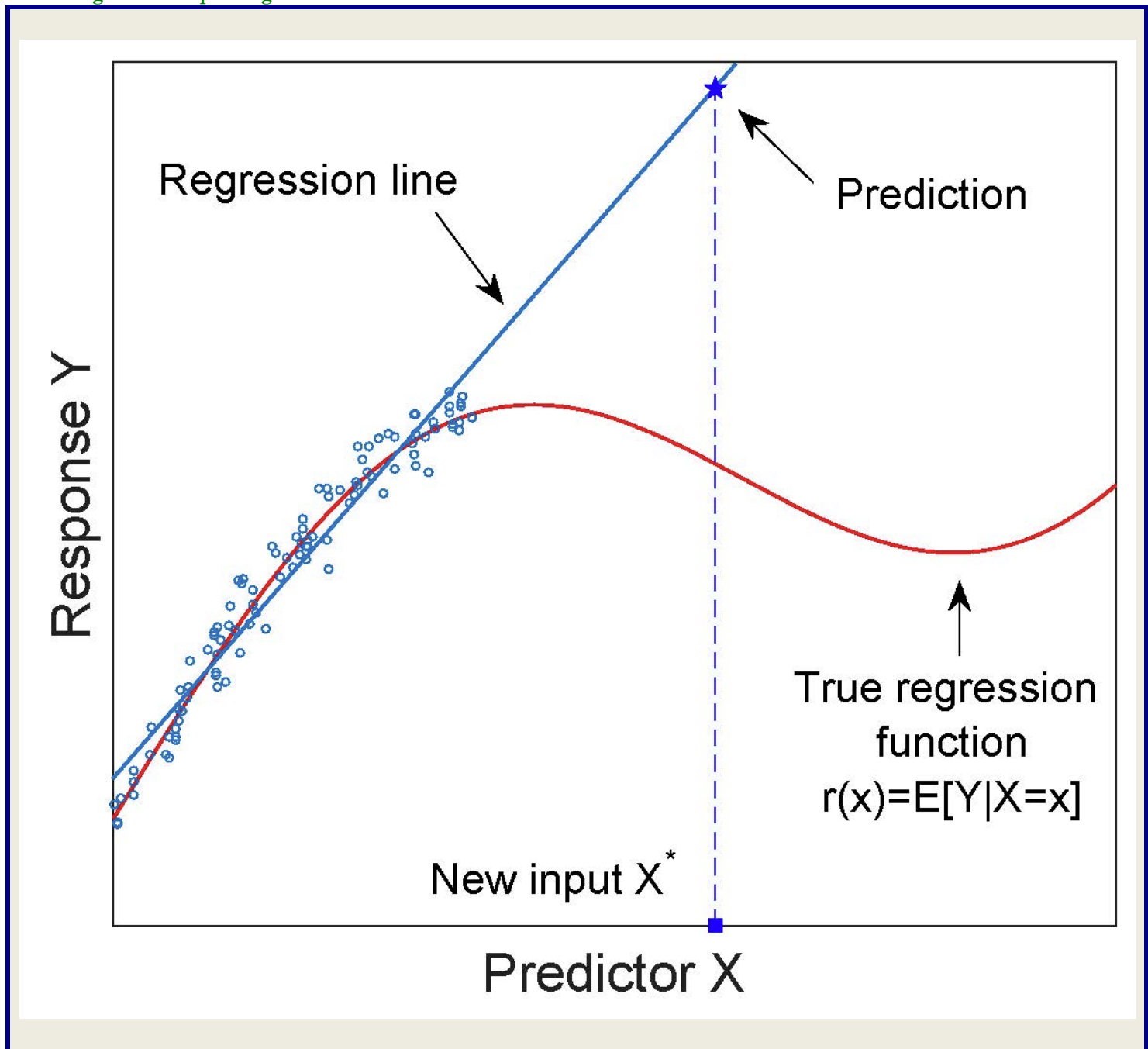
-3- **“most useful”**

Sometimes, but not so much in public health, the fitted model derives from a physical-equation. An example is Michaelis-Menton kinetics. A Michaelis-Menton model is fit to the data for the purpose of estimating the actual rate of a particular chemical reaction.

Hence – ***“A linear model is a good first start only...”***

Cheers!

The dangers of extrapolating ...

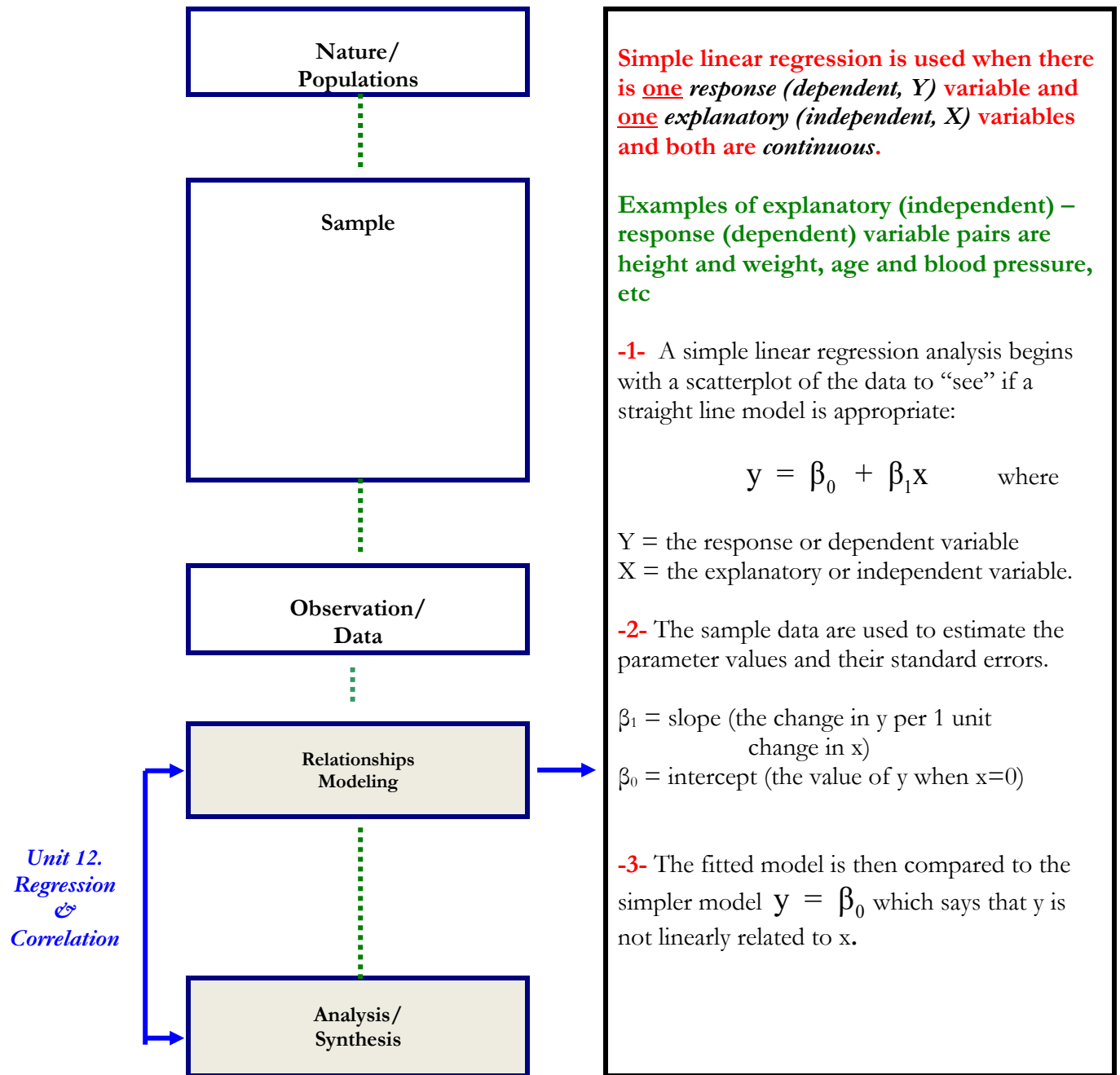


Source: Stack Exchange.

Table of Contents

		Page
Topic	1. Unit Roadmap	4
	2. Learning Objectives	5
	3. Definition of the Linear Regression Model	6
	4. Estimation	14
	5. Analysis of Variance and Introduction to R^2	26
	6. Assumptions for the Straight Line Regression	34
	7. Hypothesis Testing	38
	8. Confidence Interval Estimation	46
	9. Introduction to Correlation	51
	10. Hypothesis Test for Correlation	54

1. Unit Roadmap



2. Learning Objectives

When you have finished this unit, you should be able to:

- Explain what is meant by independent versus dependent variable and what is meant by a linear relationship;
- Produce and interpret a scatterplot;
- Define and explain the intercept and slope parameters of a linear relationship;
- Explain the theory of least squares estimation of the intercept and slope parameters of a linear relationship;
- Calculate by hand the least squares estimation of the intercept and slope parameters of a linear relationship;
- Explain the theory of the analysis of variance of simple linear regression;
- Calculate by hand the analysis of variance of simple linear regression;
- Explain, compute, and interpret R^2 in the context of simple linear regression;
- State and explain the assumptions required for estimation and hypothesis tests in regression;
- Explain, compute, and interpret the overall F-test in simple linear regression;
- Interpret the computer output of a simple linear regression analysis from a package such as R, Stata, SAS, SPSS, Minitab, etc.;
- Define and interpret the value of a Pearson Product Moment Correlation, r ;
- Explain the relationship between the Pearson product moment correlation r and the linear regression slope parameter; and
- Calculate by hand the confidence interval estimation and statistical hypothesis testing of the Pearson product moment correlation r .

3. Definition of the Linear Regression Model

Unit 11 considered two **categorical (discrete)** variables, such as smoking (yes/no) and event of low birth weight (yes/no). It was an introduction to chi-square tests of association.

Unit 12 considers two **continuous** variables, such as age and weight. It is an introduction to **simple linear regression** and **correlation**.

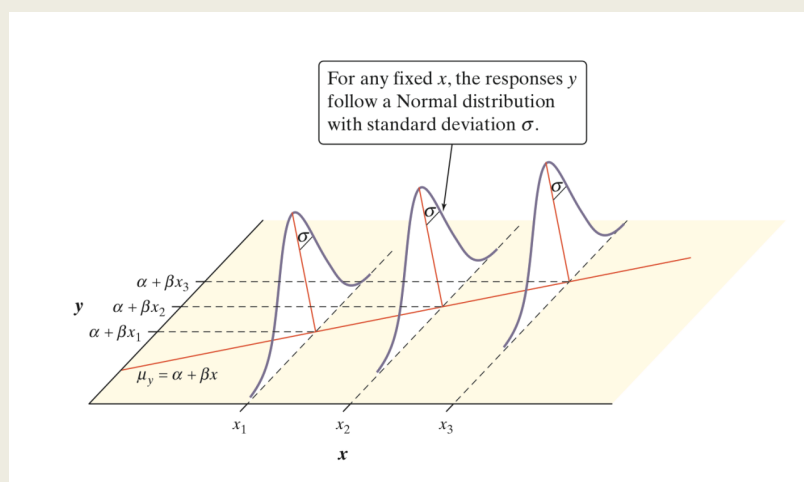
A wonderful introduction to the intuition of linear regression can be found in the text by Freedman, Pisani, and Purves (Statistics. WW Norton & Co., 1978). The following is excerpted from pp 146 and 148 of their text:

“How is weight related to height? For example, there were 411 men aged 18 to 24 in Cycle I of the Health Examination Survey. Their average height was 5 feet 8 inches = 68 inches, with an overall average weight of 158 pounds. But those men who were one inch above average in height had a somewhat higher average weight. Those men who were two inches above average in height had a still higher average weight. And so on. On the average, how much of an increase in weight is associated with each unit increase in height? The best way to get started is to look at the scattergram for these heights and weights. The object is to see how weight depends on height, so height is taken as the independent variable and plotted horizontally ...

... The regression line is to a scatter diagram as the average is to a list. The regression line estimates the average value for the dependent variable corresponding to each value of the independent variable.”

The simple linear regression model.

Consider that there is an overall distribution of Y . It has an overall mean $\mu = E[Y]$ and an overall variance $\sigma_Y^2 = \text{Var}[Y]$. Next, consider that this overall distribution is made up of subpopulations of Y , one at each level of X (for example – the distribution of Y =weight for children with X =height=50” and the distribution of Y =weight for children with X =height = 51”). We might want to know: how does the distribution of Y change, depending on which level of X we are talking about? These are called the conditional distribution of Y at X .



Source: <https://cpb-us-w2.wpmucdn.com/www.ascpages.org/dist/e/572/files/2017/03/APS-11C-12.1-LSRL-Regression-ANSWERS-1hygr7b.pdf>

Modeling the mean of Y. In simple linear regression

$$\mu_x = E[Y \text{ for the sub-population with } X = x] = E[Y | X = x] \text{ is modeled linearly in } x:$$

$$\mu_x = \beta_0 + \beta_1 x.$$

Modeling an individual observation of Y. If we have observations of Y for the subpopulation for which $X=x$, we are thus saying that each observed $Y=y$ is modeled as a departure (error) from its subpopulation-specific mean as follows:

$$y = [\text{mean}] + [\text{error in observing mean}]$$

$$= [\mu_{X=x}] + [\text{error}]$$

$$= [\beta_0 + \beta_1 x] + [\text{error}]$$

Variance of Y within each subpopulation defined by $X=x$. At each value of X, the variance of Y (we call this the conditional variance of Y) is $\sigma_{Y|X}^2$. In simple linear regression, we make the assumption that this conditional variance is the same for all subpopulations defined by X (“homogeneity of error variance”).

Correlation

Correlation considers the association of **two random** variables.

- ◆ The techniques of estimation and hypothesis testing are the same for linear regression and correlation analyses.
- ◆ Exploring the relationship begins with fitting a line to the points.

Development of a simple linear regression model analysis

Example.

Source: Kleinbaum, Kupper, and Muller 1988

The following are observations of age (days) and weight (kg) for $n=11$ chicken embryos.

WT=Y	AGE=X	LOGWT=Z
0.029	6	-1.538
0.052	7	-1.284
0.079	8	-1.102
0.125	9	-0.903
0.181	10	-0.742
0.261	11	-0.583
0.425	12	-0.372
0.738	13	-0.132
1.13	14	0.053
1.882	15	0.275
2.812	16	0.449

Nature — Population/ Sample — Observation/ Data — Relationships/ Modeling — Analysis/ Synthesis

Notation

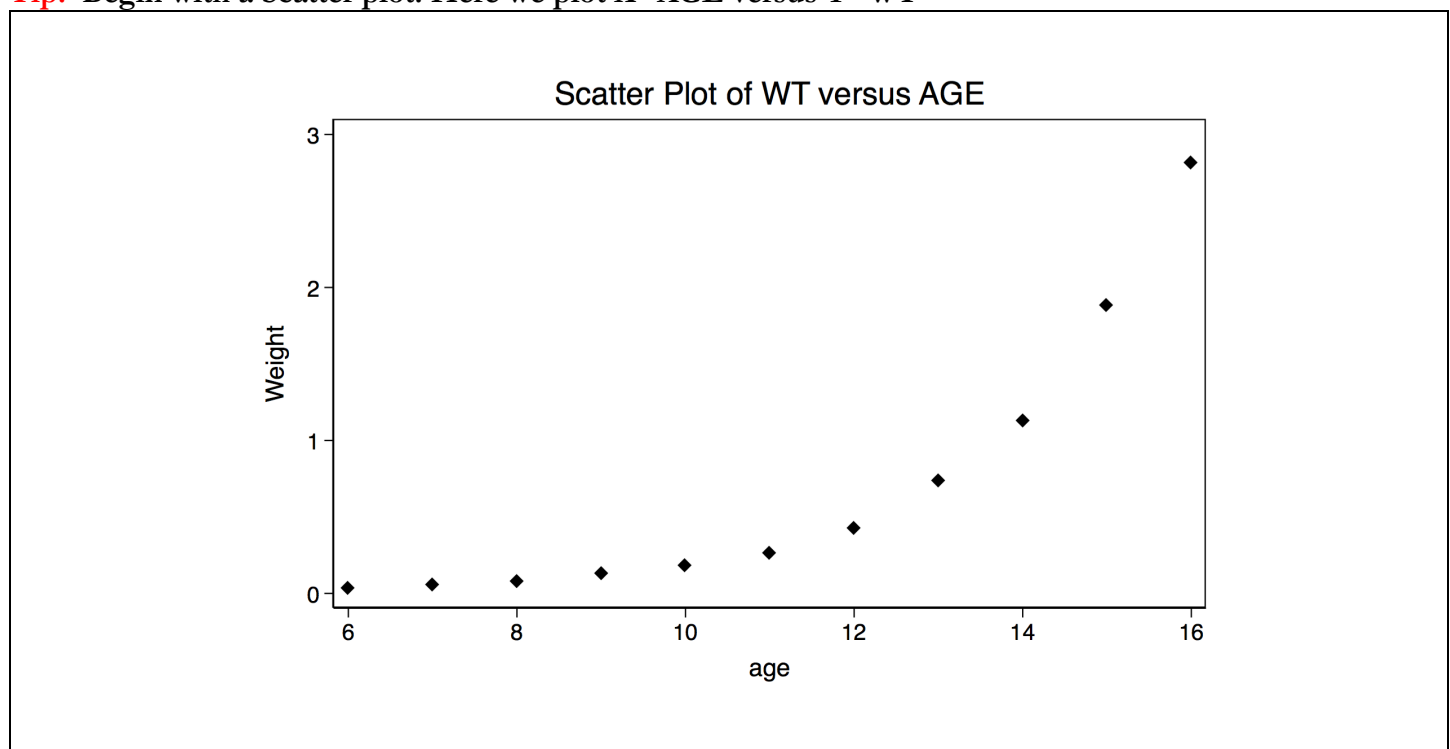
- ◆ The data are 11 pairs of (X_i, Y_i) where $X = \text{AGE}$ and $Y = \text{WT}$
 $(X_1, Y_1) = (6, .029) \cdots (X_{11}, Y_{11}) = (16, 2.812)$ and
- ◆ This table also provides 11 pairs of (X_i, Z_i) where $X = \text{AGE}$ and $Z = \text{LOGWT}$
 $(X_1, Z_1) = (6, -1.538) \cdots (X_{11}, Z_{11}) = (16, 0.449)$

Research question

There are a variety of possible research questions:

- (1) Does weight change with age?
- (2) Can the variability in weight be explained, to a significant extent, by variations in age?
- (3) What is a “good” functional form that relates age to weight?

Tip! Begin with a Scatter plot. Here we plot $X = \text{AGE}$ versus $Y = \text{WT}$



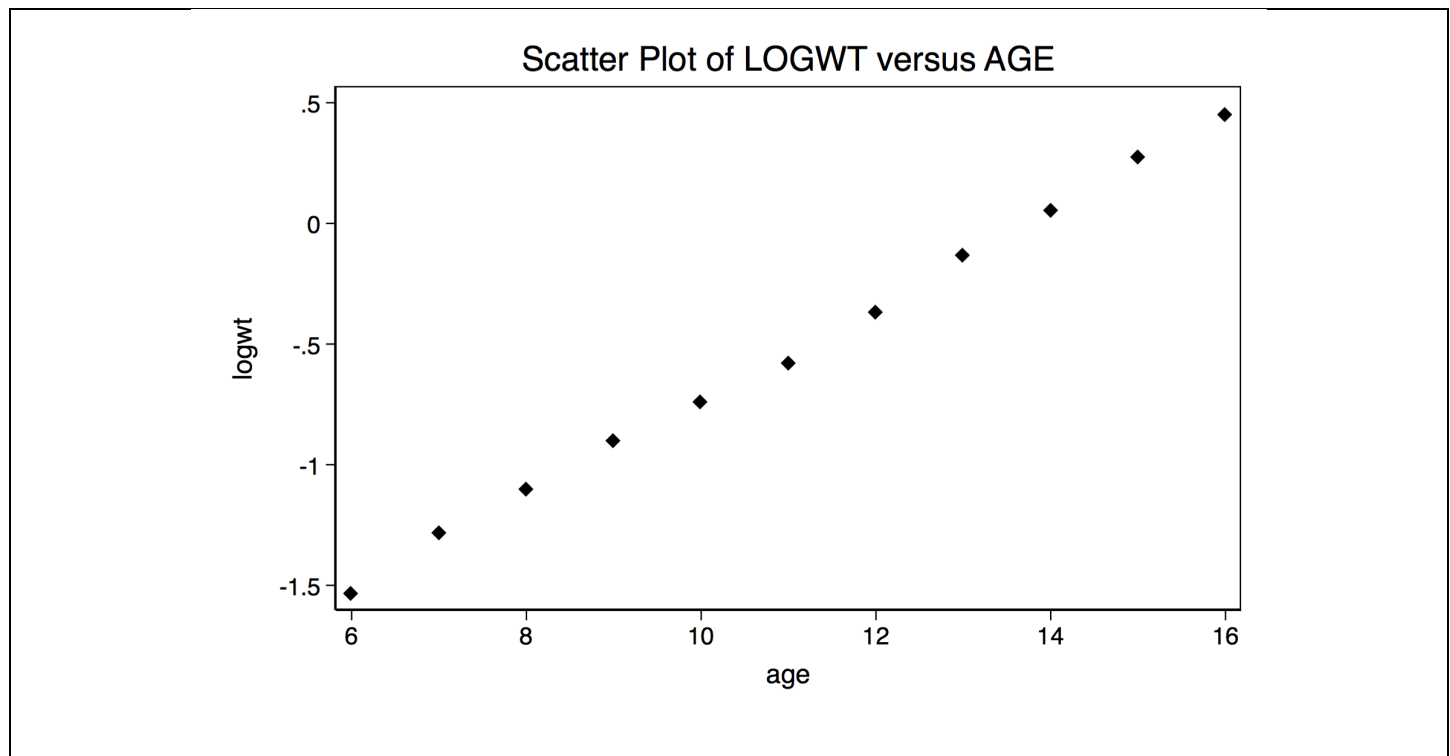
The scatterplot gives us a feel for all kinds of things!

- ◆ The average and median of X
- ◆ The range and pattern of variability in X
- ◆ The average and median of Y
- ◆ The range and pattern of variability in Y
- ◆ The nature of the relationship between X and Y
- ◆ The strength of the relationship between X and Y
- ◆ The identification of any points that might be influential

Example, continued

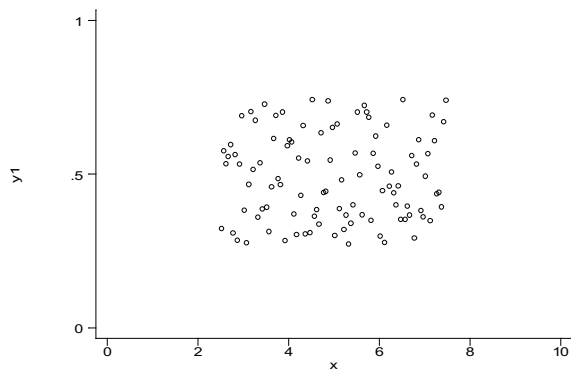
- ◆ The plot suggests a relationship between AGE and WT
- ◆ A straight line might fit well, but another model might be better
- ◆ We have adequate ranges of values for both AGE and WT
- ◆ There are no outliers

The “bowl” shape of our scatter plot suggests that perhaps a better model relates the logarithm of WT ($Z = \text{LOGWT}$) to AGE:

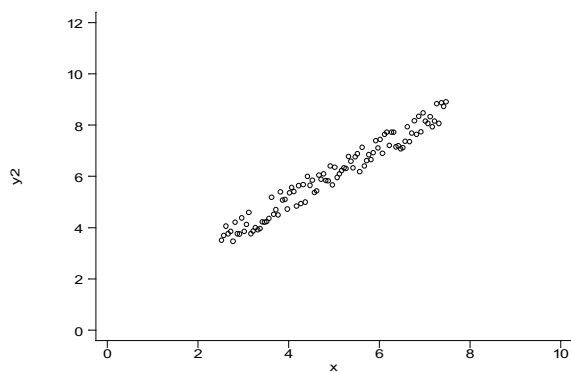


Nature
Population/
Sample
Observation/
Data
Relationships/
Modeling
Analysis/
Synthesis

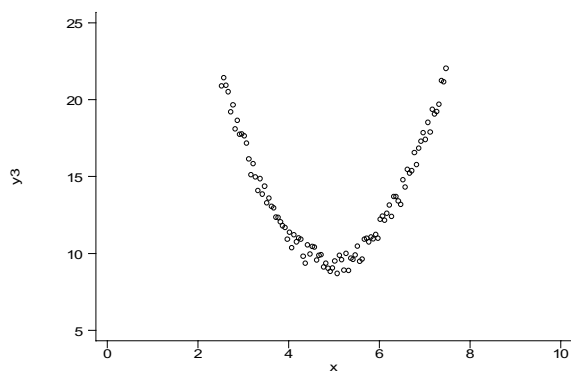
It's worth considering that we might have gotten any of a variety of plots.



No relationship between X and Y

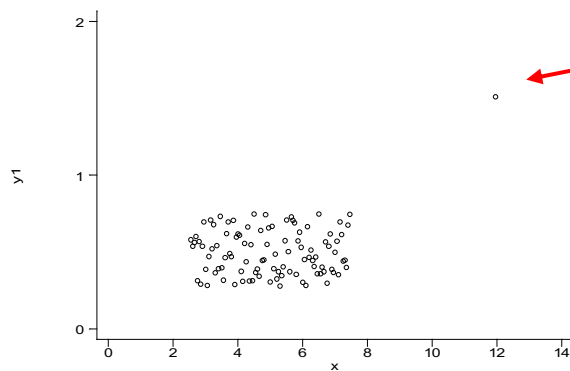


Linear relationship between X and Y



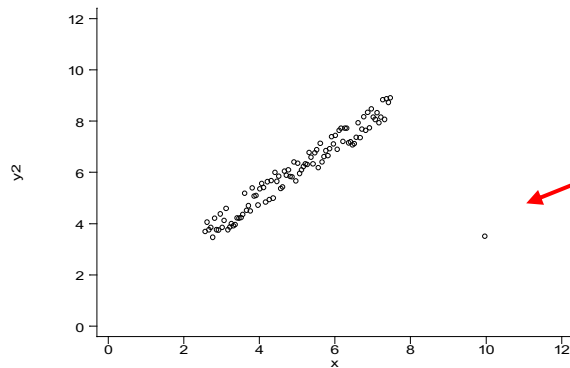
Non-linear relationship between X and Y

Nature — Population/
Sample — Observation/
Data — Relationships/
Modeling — Analysis/
Synthesis



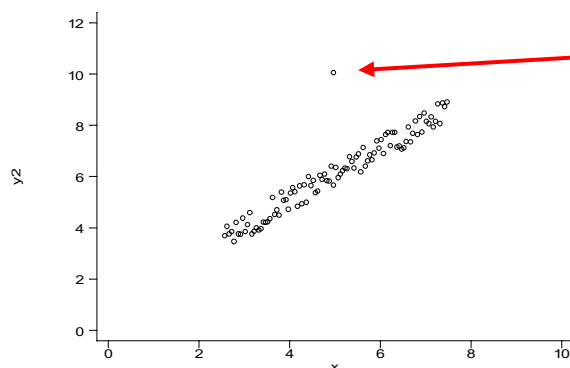
Note the outlying point

Here, a fit of a linear model will yield an estimated slope that is spuriously non-zero.



Note the outlying point

Here, a fit of a linear model will yield an estimated slope that is spuriously near zero.



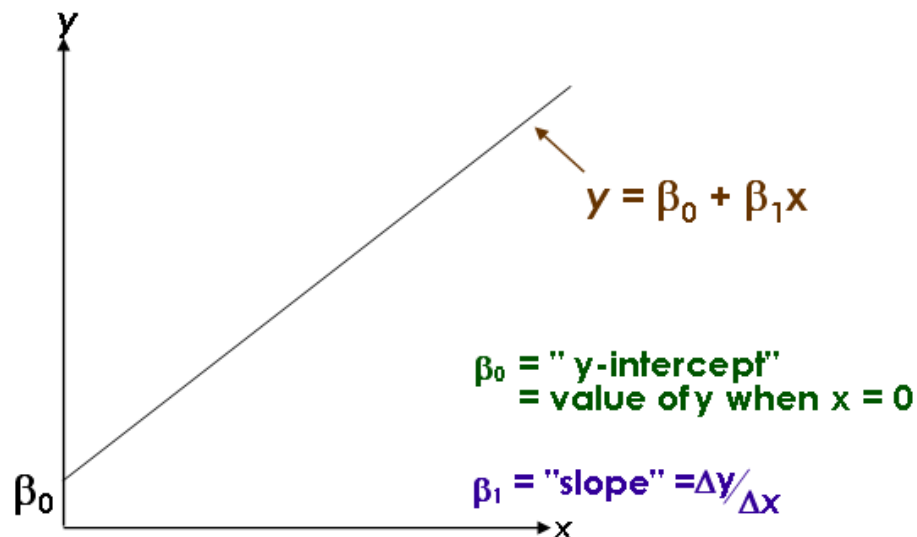
Note the outlying point

Here, a fit of a linear model will yield an estimated slope that is spuriously high.

Review of the Straight Line

Way back when, in your high school days, you may have been introduced to the straight line function, defined as “ $y = mx + b$ ” where m is the slope and b is the intercept. Nothing new here. All we’re doing is changing the notation a bit:

- (1) Slope: $m \rightarrow \beta_1$
- (2) Intercept: $b \rightarrow \beta_0$



$\beta_0 = \text{"y-intercept"} = \text{value of } y \text{ when } x = 0$

$\beta_1 = \text{"slope"} = \Delta y / \Delta x = (\text{change in } y) / (\text{change in } x)$

Slope

Slope > 0	Slope = 0	Slope < 0

Definition of the Straight Line Model

$$Y = \beta_0 + \beta_1 X$$

Population	Sample
$Y = \beta_0 + \beta_1 X + \varepsilon$	$Y = \hat{\beta}_0 + \hat{\beta}_1 X + e$
$Y = \beta_0 + \beta_1 X + \varepsilon$ = relationship in the population. $Y = \beta_0 + \beta_1 X$ is measured with <u>error ε</u> defined $\varepsilon = [Y] - [\beta_0 + \beta_1 X]$	What are $\hat{\beta}_0$, $\hat{\beta}_1$ and e ? They are our estimates of β_0 , β_1 and ε These estimates are also sometimes written as b_0 , b_1 , and e $e = \text{residual}$ is the difference between the observed and the estimated model $e = [Y] - [\hat{\beta}_0 + \hat{\beta}_1 X]$
β_0 , β_1 and ε are all <u>unknown!!</u>	We obtain the estimates $\hat{\beta}_0$, $\hat{\beta}_1$ and e by the method of <u>least squares estimation</u> .
	$\hat{\beta}_0$, $\hat{\beta}_1$ and e are <u>known</u> How close did we get? To see if $\hat{\beta}_0 \approx \beta_0$ and $\hat{\beta}_1 \approx \beta_1$ we perform <u>regression diagnostics</u> . <i>Regression diagnostics are discussed in BIOSTATS 640</i>

Notation ... sorry ...

Y = the outcome or dependent variable

X = the predictor or independent variable

μ_Y = The expected value of Y for all persons in the population

$\mu_{Y|X=x}$ = The expected value of Y for the sub-population for whom $X=x$

σ_Y^2 = Variability of Y among all persons in the population

$\sigma_{Y|X=x}^2$ = Variability of Y for the sub-population for whom $X=x$

Nature Population/
Sample Observation/
Data Relationships/
Modeling Analysis/
Synthesis

4. Estimation

Least squares estimation is used to obtain guesses of β_0 and β_1 .

When the outcome = Y is distributed normal, least squares estimation is the same as maximum likelihood estimation. **Note – If you are not familiar with “maximum likelihood estimation”, don’t worry. This is introduced in BIOSTATS 640.**

“Least Squares”, “Close” and Least Squares Estimation

Theoretically, it is possible to draw many lines through an X-Y scatter of points. Which to choose? “Least squares” estimation is one approach (fyi – there are others) to choosing a line that is a good fit to the data.

- ♦ $d_i = [\text{observed } Y - \text{fitted } \hat{Y}]$ for the i^{th} person
Perhaps we’d like $d_i = [\text{observed } Y - \text{fitted } \hat{Y}] = \text{smallest possible}$.
Note that this is a vertical distance, since it is a distance on the vertical axis.
- ♦ $d_i^2 = [Y_i - \hat{Y}_i]^2$
Better yet, perhaps we’d like to minimize the squared difference:
 $d_i^2 = [\text{observed } Y - \text{fitted } \hat{Y}]^2 = \text{smallest possible}$
- ♦ **Glitch.** We can’t minimize each d_i^2 separately. In particular, it is not possible to choose common values of $\hat{\beta}_0$ and $\hat{\beta}_1$ that minimizes

$$d_1^2 = (Y_1 - \hat{Y}_1)^2 \quad \text{for subject 1 \textit{and} minimizes}$$

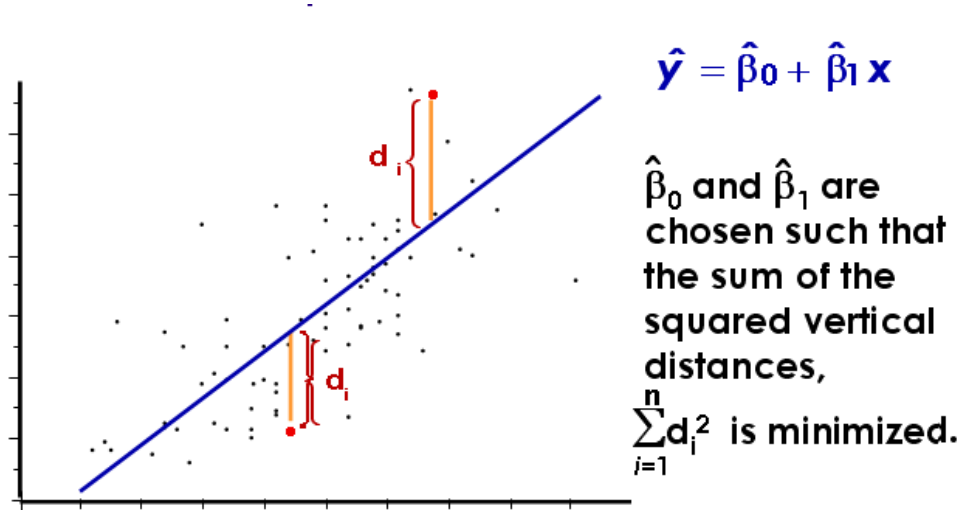
$$d_2^2 = (Y_2 - \hat{Y}_2)^2 \quad \text{for subject 2 \textit{and} minimizes}$$

$$\dots \dots \textit{and} \text{ minimizes}$$

$$d_n^2 = (Y_n - \hat{Y}_n)^2 \quad \text{for the } n^{\text{th}} \text{ subject}$$

- ♦ So, instead, we choose values for $\hat{\beta}_0$ and $\hat{\beta}_1$ that, upon insertion, minimizes the total

$$\sum_{i=1}^n d_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - [\hat{\beta}_0 + \hat{\beta}_1 X_i])^2$$



For each observed value x_i , we have an observed y_i , and the “predicted” value \hat{y}_i , on the line. The vertical distances $d_i = (y_i - \hat{y}_i)$.

$$\sum_{i=1}^n d_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n \left(Y_i - [\hat{\beta}_0 + \hat{\beta}_1 X_i] \right)^2 \text{ has a variety of names:}$$

- ◆ residual sum of squares, SSE or SSQ(residual)
- ◆ sum of squares about the regression line
- ◆ sum of squares due error (SSE)

A closer look ...

Some very helpful preliminary calculations

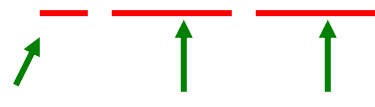
- $S_{xx} = \sum (X - \bar{X})^2 = \sum X^2 - N\bar{X}^2$
- $S_{yy} = \sum (Y - \bar{Y})^2 = \sum Y^2 - N\bar{Y}^2$
- $S_{xy} = \sum (X - \bar{X})(Y - \bar{Y}) = \sum XY - N\bar{X}\bar{Y}$

Note - These expressions make use of a “summation notation”, introduced in Unit 1.

*The capital “S” indicates “**summation**”.*

*In S_{xy} , the first subscript “**x**” is saying $(x - \bar{x})$.*

*The second subscript “**y**” is saying $(y - \bar{y})$.*

$$S_{xy} = \sum (X - \bar{X})(Y - \bar{Y})$$


S subscript x subscript y

Estimate of Slope	$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\text{cov}(X, Y)}{\text{var}(X)}$	$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$
Estimate of Intercept	$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$	
Prediction of Y	$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$ $= b_0 + b_1 X$	

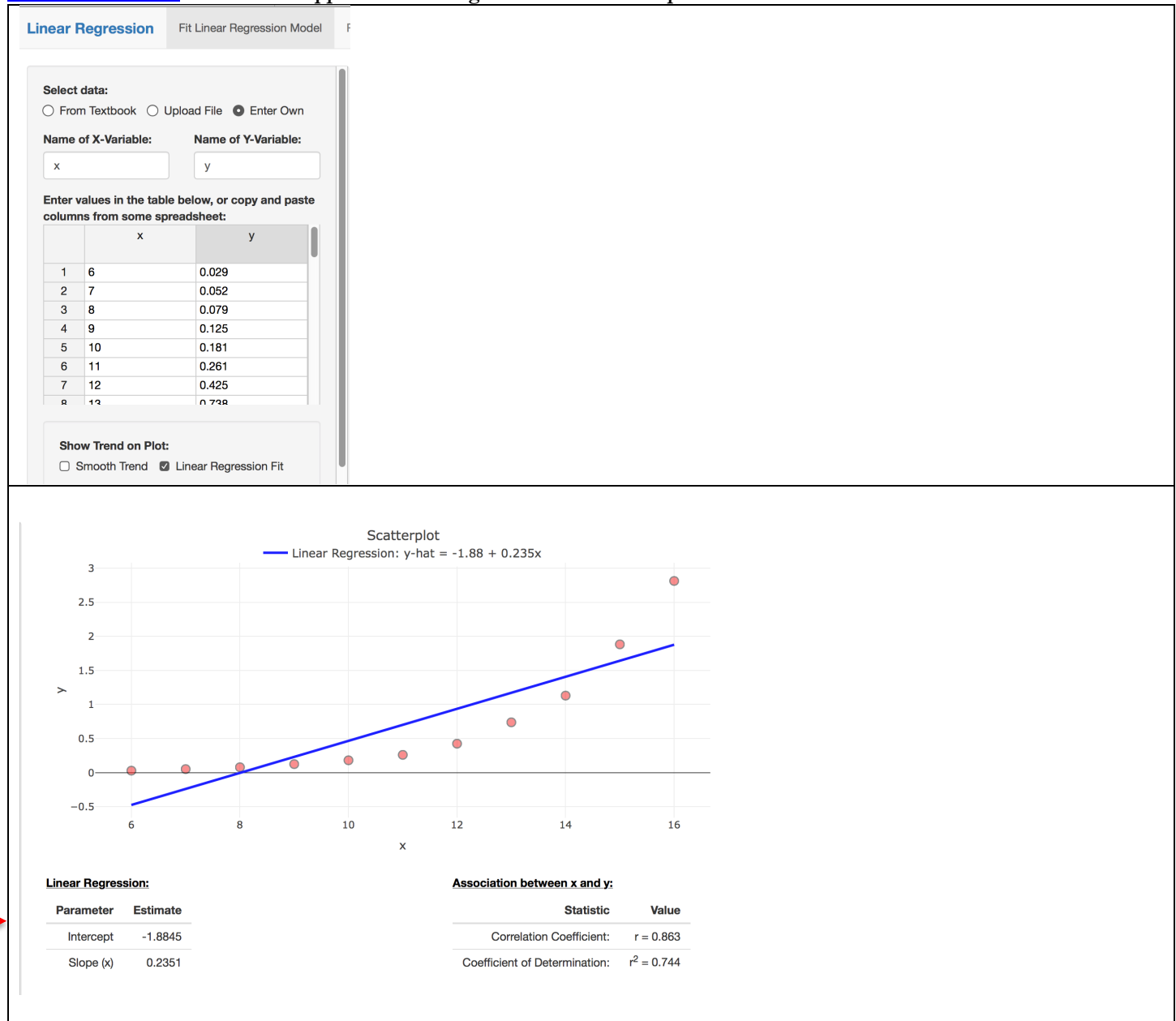
Do these estimates make sense?

Estimate of Slope	$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\text{cov}(X, Y)}{\text{var}(X)}$	<p>The linear movement in Y with linear movement in X is measured relative to the variability in X.</p> <p>$\hat{\beta}_1 = 0$ says: With a unit change in X, overall there is a 50-50 chance that Y increases versus decreases</p> <p>$\hat{\beta}_1 \neq 0$ says: With a unit increase in X, Y increases also ($\hat{\beta}_1 > 0$) or Y decreases ($\hat{\beta}_1 < 0$).</p>
Estimate of Intercept	$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$	<p>If the linear model is incorrect, or, if the true model does not have a linear component, we obtain $\hat{\beta}_1 = 0$ and $\hat{\beta}_0 = \bar{Y}$ as our best guess of an unknown Y</p>

ILLUSTRATION of Model Estimation: $Y=WT$ and $X=AGE$

Art of Stat

www.artofstat.com > online webapps > Linear Regression > At left drop down Enter Data: "Enter Own"



The fitted line is therefore

$$y = -1.8845 + 0.2351 \cdot x.$$

Key: For each ONE unit (1 day) increase in x, y is estimated to increase by 0.2351 units (kg)

R

```
setwd("/Users/cbigelow/Desktop/") # Set working directory (yours will be different)
rm(list=ls()) # Clear the decks (clear the working environment)
options(scipen=1000) # Turn off scientific notation
options(show.signif.stars=FALSE) # Turn off display of significance stars
```

Input data: Copy/paste from Excel -> table -> data frame

```
datatable=read.table(text="
y_wt    x_age    z_logwt
0.029    6.000    -1.538
0.052    7.000    -1.284
0.079    8.000    -1.102
0.125    9.000    -0.903
0.181    10.000   -0.742
0.261    11.000   -0.583
0.425    12.000   -0.372
0.738    13.000   -0.132
1.130    14.000    0.053
1.882    15.000    0.275
2.812    16.000    0.449",header=TRUE)
dataset <- as.data.frame.matrix(datatable)
```

Fit Simple Linear Regression: Dependent=y_wt Predictor=x_age

```
fit1 <- lm(y_wt ~ x_age, data=dataset) # lm( ) will fit the model. Result is stored in object named fit1
summary(fit1)

##
## Call:
## lm(formula = y_wt ~ x_age, data = dataset)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.5113 -0.3593 -0.1061  0.2657  0.9354
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.88453    0.52584  -3.584 0.005895
## x_age        0.23507    0.04594   5.117 0.000631
##
## Residual standard error: 0.4818 on 9 degrees of freedom
## Multiple R-squared:  0.7442, Adjusted R-squared:  0.7158
## F-statistic: 26.18 on 1 and 9 DF,  p-value: 0.0006308
```

Here (similar to the artofstat output), the fitted line is

$$y_wt = -1.8845 + 0.2351 \cdot x_age.$$

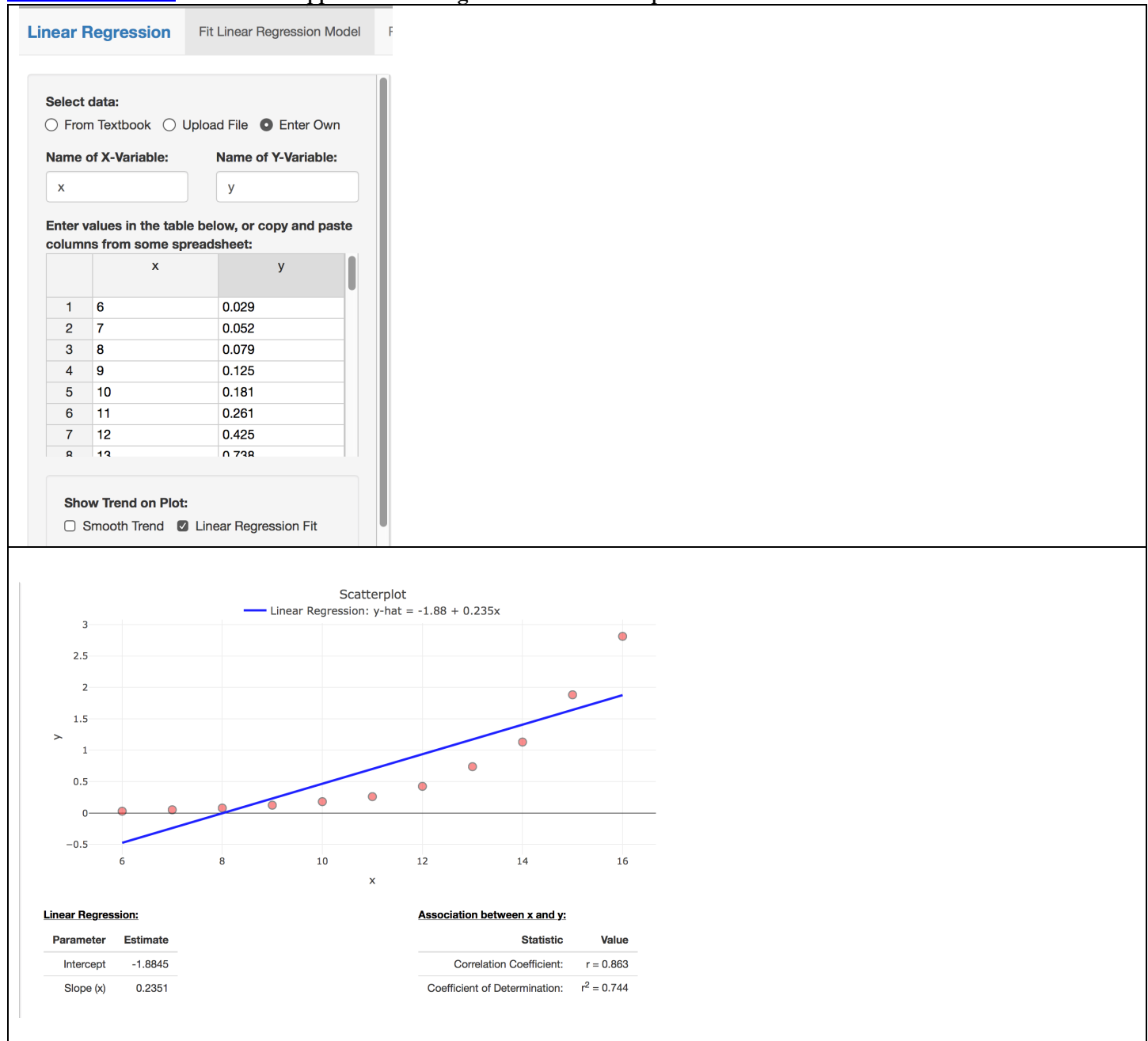
KEY:: For each ONE unit (1 day) increase in x_age, y_wt is estimated to increase by 0.2351 units (kg)



ILLUSTRATION of Plot of Scatter with Overlay Fit: $Y=WT$ and $X=AGE$

Art of Statt

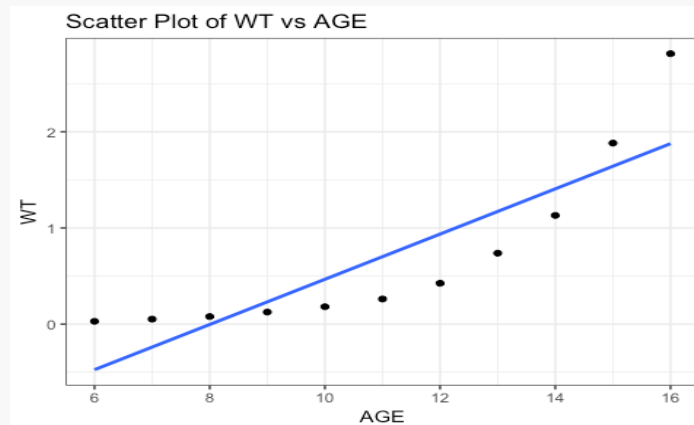
www.artofstat.com > online webapps > Linear Regression > At left drop down Enter Data: "Enter Own"



Nature _____ Population/ Sample _____ Observation/ Data _____ Relationships/ Modeling _____ Analysis/ Synthesis

R

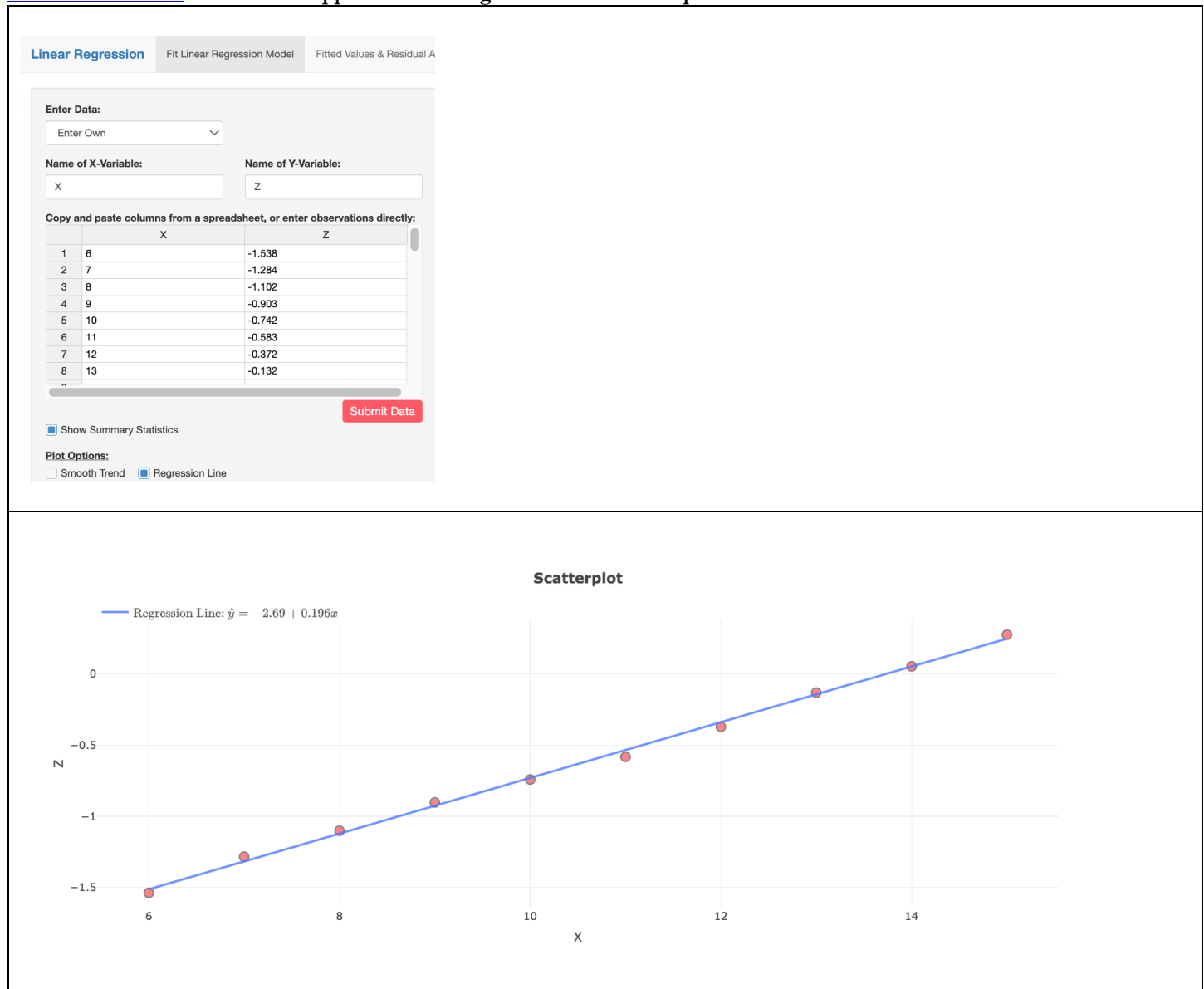
```
library(ggplot2) # Library( ) to attach the package {ggplot2} to session (one time)
ggplot(dataset, aes(x=x_age,y=y_wt)) +
  geom_smooth(method=lm, se=FALSE) +
  geom_point() +
  xlab("AGE") +
  ylab("WT") +
  ggtitle("Scatter Plot of WT vs AGE") +
  theme_bw()
```



- ◆ As we might have guessed, the straight-line model may not be the best choice.
- ◆ The “bowl” shape of the scatter plot does have a linear component, however.
- ◆ Without the plot, we might have believed the straight-line fit is okay.

Art of Stat. Z=LOGWT and X=AGE

www.artofstat.com > online webapps > Linear Regression > At left drop down Enter Data: "Enter Own"



The fitted line is

$$Z = -2.69 + 0.196 \cdot x$$

Key: For each ONE unit (1 day) increase in x, z is estimated to increase by 0.1959

Nature
Population/
Sample
Observation/
Data
Relationships/
Modeling
Analysis/
Synthesis

R. Z=LOGWT and X=AGE

```
# Fit Simple Linear Regression: Dependent=z_logwt Predictor=x_age

fit2 <- lm(z_logwt ~ x_age, data=dataset)
summary(fit2)

##
## Call:
## lm(formula = z_logwt ~ x_age, data = dataset)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.04854 -0.01787  0.00400  0.02168  0.03402
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept) -2.689255    0.030637  -87.78 0.0000000000000164
## x_age        0.195891    0.002677   73.18 0.0000000000000840
##
## Residual standard error: 0.02807 on 9 degrees of freedom
## Multiple R-squared:  0.9983, Adjusted R-squared:  0.9981
## F-statistic: 5356 on 1 and 9 DF, p-value: 0.0000000000008399
```

The fitted line is

$$z_logwt = -2.6892 + 0.1959 \cdot x_age.$$

Key: For each ONE unit (1 day) increase in x_age , z_logwt is estimated to increase by 0.1959

R ILLUSTRATION of Plot of Scatter with Overlay Fit: Z=LOGWT and X=AGE

```
library(ggplot2)
ggplot(dataset, aes(x=x_age,y=z_logwt)) +
  geom_smooth(method=lm, se=FALSE) +
  geom_point() +
  xlab("AGE") +
  ylab("LOGWT") +
  ggtitle("Scatter Plot of LOGWT vs AGE") +
  theme_bw()
```

geom_smooth(method=lm, se=FALSE) plots fitted line
geom_point() produces x-y scatter



For the brave – Try doing the calculations by hand ...

Prediction of Weight from Height

Source: Dixon and Massey (1969)

Individual	Height (X)	Weight (Y)
1	60	110
2	60	135
3	60	120
4	62	120
5	62	140
6	62	130
7	62	135
8	64	150
9	64	145
10	70	170
11	70	185
12	70	160

Preliminary calculations

$\bar{X} = 63.833$	$\bar{Y} = 141.667$
$\sum X_i^2 = 49,068$	$\sum Y_i^2 = 246,100$
$\sum X_i Y_i = 109,380$	$S_{xx} = 171.667$
$S_{yy} = 5,266.667$	$S_{xy} = 863.333$

Estimate of Slope	$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$	$\hat{\beta}_1 = \frac{863.333}{171.667} = 5.0291$
Estimate of Intercept	$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$	$\hat{\beta}_0 = 141.667 - (5.0291)(63.8333) = -179.3573$

Nature _____ Population/ Sample _____ Observation/ Data _____ Relationships/ Modeling _____ Analysis/ Synthesis

5. Analysis of Variance and Introduction to R^2

Analysis of Variance

One goal (by no means the only one!) is to explain the variability in our outcomes Y .

The outcomes are the values of the dependent variable Y . In the example on page 7, the outcomes were $y_1 = 0.029, y_2 = 0.052, \dots, y_{11} = 2.812$. In fitting a simple linear regression of these weights in the predictor $X = \text{age}$, our goal (*one of them*) was to learn if some of the variability in weights could be explained by age.

The variability in our outcomes Y that we seek to explain is called the “total sum of squares in Y ”, also called the “total sum of squares, corrected”.

Total Variability “to be explained”
Total Sum of Squares

$$\sum_{i=1}^n (y_i - \bar{y})^2$$

Key – this is the total variability of the individual observed y about their average \bar{y}

- Features
 - Notice that the Total Sum of Squares = the numerator of the sample variance of the Y 's
 - Because there's no division by anything, you can think of the total sum of squares as a measure of *total* scatter
 - Another way of thinking of it is to think of Total Sum of Squares = total “noisiness” of the outcomes y_1, y_2, \dots, y_n
- The total sum of squares goes by several names and notations (sorry!)
 - “Total sum of squares”
 - “Total sum of squares, corrected”
 - SSY
 - SST

Nature _____ Population/
Sample _____ Observation/
Data _____ Relationships/
Modeling _____ Analysis/
Synthesis

The analysis of variance starts with this total sum of squares = $\sum_{i=1}^n (Y_i - \bar{Y})^2$ and, like a (delicious) pie, partitions (carves) it into components (wedges).

In simple linear regression, the total is partitioned into just 2 components (wedges of the pie):

1. **Due residual** (the individual Y about the individual prediction \hat{Y})
2. **Due regression** (the prediction \hat{Y} about the overall mean \bar{Y})

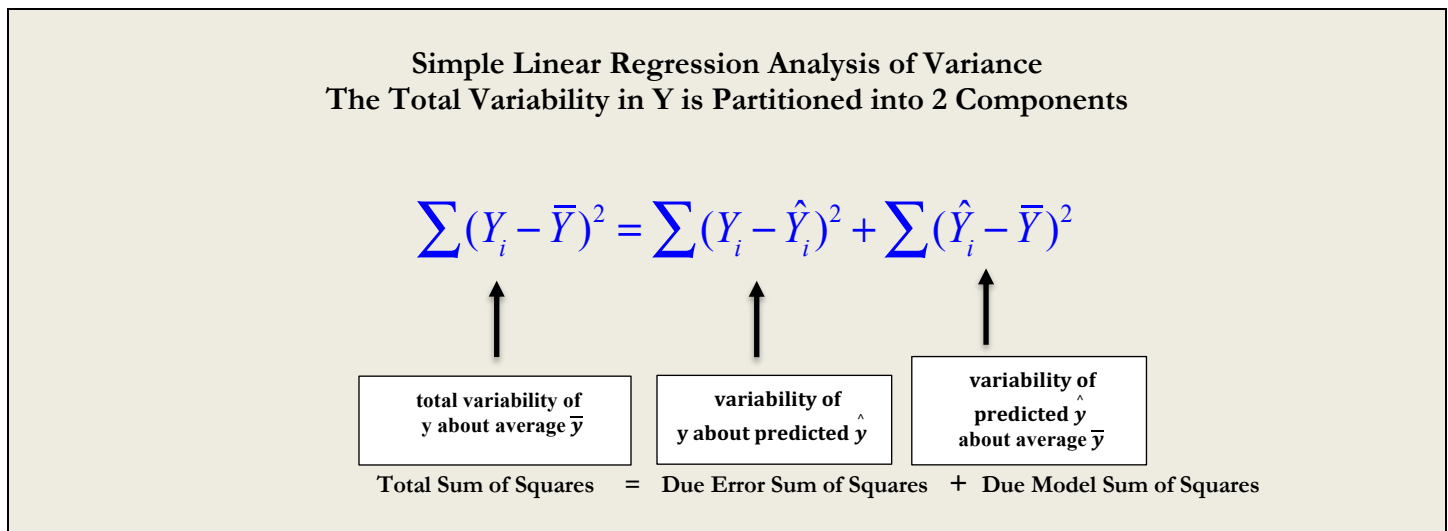
Here is the partition (Note – Look closely and you’ll see that both sides are the same)

$$(Y_i - \bar{Y}) = (Y_i - \hat{Y}_i) + (\hat{Y}_i - \bar{Y})$$

Some algebra (not shown) reveals a nice partition of the total variability.

$$\sum (Y_i - \bar{Y})^2 = \sum (Y_i - \hat{Y}_i)^2 + \sum (\hat{Y}_i - \bar{Y})^2$$

Total Sum of Squares = Due Error Sum of Squares + Due Model Sum of Squares



Example.

Consider again the data on page 7 and, specifically, the simple linear model where outcomes $Y=WT$ were modeled linearly in $X=AGE$.

Art of Stat. $Y=WT$ and $X=AGE$

www.artofstat.com > online webapps > Linear Regression > At left drop down Enter Data: "Enter Own"

Linear Regression | Fit Linear Regression Model | Fitted Values & Residual A

Enter Data:
Enter Own

Name of X-Variable: X Name of Y-Variable: Y

Copy and paste columns from a spreadsheet, or enter observations directly:

	X	Y
1	6	0.029
2	7	0.052
3	8	0.079
4	9	0.125
5	10	0.181
6	11	0.261
7	12	0.425
8	13	0.738

Linear Regression Equation:

Parameter	Estimate
Intercept	-1.885
Slope (X)	0.2351

Model Summary:

Statistic	Value
Correlation Coefficient r	0.863
Coefficient of Determination r^2	74.4%
Residual Standard Deviation	0.4818

ANOVA Table:

Source	Df	Sum Sq	Mean Sq	F value	P-value
Regression	1	6.1	6.08	26.2	0.00
Residuals	9	2.1	0.23		
Total	10	8.2			

Regression Options:

☐ Find Predicted Value

☐ Show Residuals on Plot

☐ Show Standard Errors & P-values

☐ Confidence Interval for Slope

☐ Confidence/Prediction Interval

☒ ANOVA Table

R

Note: The command `anova()` does not produce display of the total sum of squares.

```

fit1 <- lm(y_wt ~ x_age, data=dataset)
summary(fit1)
anova(fit1)

```

Analysis of Variance Table

Response: y_wt

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
x_age	1	6.0785	6.0785	26.18	0.0006308
Residuals	9	2.0896	0.2322		

KEY:

Due Model Sum of Squares = $\sum(\hat{y} - \bar{y})^2$ = x_age Sum Sq = 6.0785

Due Error Sum of Squares = $\sum(y - \hat{y})^2$ = Residuals Sum Sq = 2.0896

A closer look...

Total Sum of Squares = Due Model Sum of Squares + Due Error Sum of Squares

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Total
due model
due error
sum of squares
sum of squares
sum of squares

- ◆ $(Y_i - \bar{Y})$ = deviation of Y_i from \bar{Y} that is to be explained
- ◆ $(\hat{Y}_i - \bar{Y})$ = “due model”, “signal”, “systematic”, “due regression”
- ◆ $(Y_i - \hat{Y}_i)$ = “due error”, “noise”, or “residual”

We hope that we can **explain** a lot of the total variability $\sum_{i=1}^n (Y_i - \bar{Y})^2$ with a fitted model:

What happens when $\beta_1 \neq 0$?	What happens when $\beta_1 = 0$?
We have an actual slope!	The slope is zero.
A straight-line relationship is helpful	A straight-line relationship is not helpful
Best guess is $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$	Best guess is $\hat{Y} = \hat{\beta}_0 = \bar{Y}$
Due model “sum of squares” tends to be LARGE because $(\hat{Y} - \bar{Y}) = (\hat{\beta}_0 + \hat{\beta}_1 X) - \bar{Y}$ $= \bar{Y} - \hat{\beta}_1 \bar{X} + \hat{\beta}_1 X - \bar{Y}$ $= \hat{\beta}_1 (X - \bar{X})$	Due error “sum of squares” tends to be nearly the TOTAL because $(Y - \hat{Y}) = (Y - [\hat{\beta}_0]) = (Y - \bar{Y})$
Due error “sum of squares” has to be small	Due regression “sum of squares” has to be small
\rightarrow $\frac{\text{due(model)}}{\text{due(error)}}$ will be large	\rightarrow $\frac{\text{due(model)}}{\text{due(error)}}$ will be small

How to Partition the Total Variance into
[what the model explains] + [what is left-over, unexplained]
all things sums of squares and mean squares

1. Total Variance. *This is the total sum of squares. Think of it as the “whole pie”*

- ◆ $\sum_{i=1}^n (Y_i - \bar{Y})^2 = \text{“total sum of squares(corrected)”}$
- ◆ Degrees of freedom = df = (n-1)
- ◆ $\sum_{i=1}^n (Y_i - \bar{Y})^2 / [n - 1] = \text{“total mean square”}$ Note – This is NOT displayed in the analysis of variance table

2. Due Model. *This is what the model explains. Think of it as “one piece of the pie”.*

- ◆ $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = \hat{\beta}_1^2 \sum_{i=1}^n (X_i - \bar{X})^2 = \text{“regression sum of squares”}$
- ◆ Degrees of freedom = df = 1
- ◆ $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = \hat{\beta}_1^2 \sum_{i=1}^n (X_i - \bar{X})^2 / [1] = \text{“regression mean square”}$, also called the
“model mean square”. It is an example of a variance component.

3. Due Residual/Error. *This is what is left-over. It is what the model does NOT explain. Think of it as the “remaining/other piece of the pie”*

- ◆ $\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \text{“residual sum of squares”}$
- ◆ Degrees of freedom = df = (n-2)
- ◆ $\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 / [n - 2] = \text{“residual mean square”}$, also called the “error mean square”

Source	df	Sum of Squares A measure of variability	Mean Square = Sum of Squares / df A measure of average/typical/mean variability
Regression due model	1	$SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$	$msq(model) = SSR/1$
Residual due error	(n-2)	$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$	$msq(residual) = SSE/(n-2) = \hat{\sigma}_{Y X}^2$
Total, corrected	(n-1)	$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$	

Be careful! The question we may ask from an analysis of variance table is a limited one.

Does the fit of the straight-line model explain a significant portion of the variability of the individual Y about \bar{Y} ?

Is this fitted model better than using \bar{Y} alone?

We are NOT asking:

Is the choice of the straight line model correct? *nor are we asking*
Would another functional form be a better choice?

We'll use a hypothesis test approach (another “proof by contradiction” reasoning just like we did in Units 8-10).

- ◆ Step 1: Assume, provisionally, the “nothing is going on” null hypothesis that says $\beta_1 = 0$ (“no linear relationship”)
- ◆ Step 2: Use least squares estimation to estimate a “closest” line
- ◆ Step 3: Do the partition of the total sum of squares and produce the analysis of variance table. Using the analysis of variance table, compare the due regression mean square to the residual mean square
- ◆ Step 4: The null hypothesis says the slope is zero. Consider what happens vis a vis the slope β_1 ?

Null true $\rightarrow \beta_1 = 0 \rightarrow$ the due (regression)/due (residual) will be SMALL

Null NOT true $\rightarrow \beta_1 \neq 0 \rightarrow$ the due (regression)/due (residual) will be LARGE

- ◆ Step 5: Carry out a p-value calculation that assumes the null is true and that answers the following question (just as we did in Units 8, 9 and 10)

If the null hypothesis is true and $\beta_1 = 0$ truly, what were the chances of obtaining a value of due (regression)/due (residual) as large or larger than that observed?

*To calculate “chances of extremeness under some assumed null hypothesis”
we need a null hypothesis probability model!
But did you notice? So far, we have not actually used one!*

Nature — Population/
Sample — Observation/
Data — Relationships/
Modeling — Analysis/
Synthesis

R²

R² is the proportion (%) of the total variability in Y that is explained by the fitted model

R²
Coefficient of Determination

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\text{Due Model Sum of Squares}}{\text{Total Sum of Squares}} = \frac{SSR}{SST}$$

Note – This is often multiplied by 100 and expressed as a %

- Features
 - R² is the percent of the total variability that is explained by the model just fit
 - It is a proportion
- Special Case: Simple Linear Regression
 - $\sqrt{R^2} = r$ = Pearson product moment correlation
 - r is a measure of linear association
 - More on this ahead in Section 9. *Introduction to Correlation*

Art of Stat. Y=WT and X=AGE

www.artofstat.com > online webapps > Linear Regression > At left drop down Enter Data: “Enter Own”

Model Summary:

Statistic	Value
Correlation Coefficient <i>r</i>	0.863
Coefficient of Determination <i>r</i> ²	74.4%
Residual Standard Deviation	0.4818



Nature _____
 Population/
Sample _____
 Observation/
Data _____
 Relationships/
Modeling _____
 Analysis/
Synthesis

R

```
fit1 <- lm(y_wt ~ x_age, data=dataset)
summary(fit1)
```

```
##
## Residual standard error: 0.4818 on 9 degrees of freedom
## Multiple R-squared:  0.7442, Adjusted R-squared:  0.7158
## F-statistic: 26.18 on 1 and 9 DF, p-value: 0.0006308
```

KEY:

Due Model Sum of Squares = $\sum (\hat{y} - \bar{y})^2$ = x_age Sum Sq = 6.0785

Due Error Sum of Squares = $\sum (y - \hat{y})^2$ = Residuals Sum Sq = 2.0896

R Squared = [Model Sum of Squares] / [Total Sum of Squares] = 6.0785 / [6.0785 + 2.0896] = 0.7442

6. Assumptions for a Straight-Line Regression Analysis

In doing least squares estimation, we did not use a probability model. We were minimizing vertical distances (geometry). If we want to do confidence interval estimation and/or test some null hypotheses, we need to have a probability model and some assumptions. Here you go!

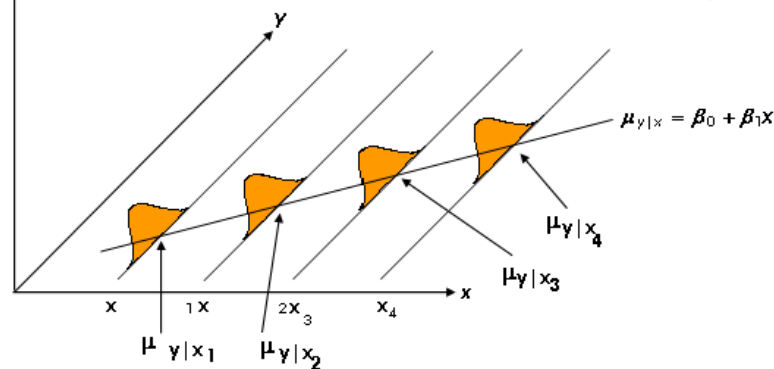
Assumptions for Simple Linear Regression

- ◆ The separate observations Y_1, Y_2, \dots, Y_n are independent.
- ◆ The values of the predictor variable X are fixed and measured without error.
- ◆ For each value of the predictor variable $X=x$, the distribution of values of Y follows a normal distribution with mean equal to $\mu_{Y|X=x}$ and common variance equal to $\sigma_{Y|x}^2$.
- ◆ The separate means $\mu_{Y|X=x}$ lie on a straight line; that is –

$$\mu_{Y|X=x} = \beta_0 + \beta_1 X$$

At each value of X , there is a population of Y for persons with $X=x$

For each value of x , the values of y are normally distributed around $\mu_{Y|x}$, on the line, with the same variance for all values of x , but different means, $\mu_{Y|x}$.



Here, $\sigma_{Y|x_1}^2 = \sigma_{Y|x_2}^2 = \sigma_{Y|x_3}^2 = \sigma_{Y|x_4}^2$

With these assumptions, we can assess the significance of the variance explained by the model.

$$F = \frac{\text{mean square(model)}}{\text{mean square(residual)}} = \frac{\text{msq(model)}}{\text{msq(residual)}} \quad \text{with df} = 1, (n-2)$$

When $\beta_1 = 0$ The slope is zero, meaning there is NO linearity	When $\beta_1 \neq 0$ We have a slope! This means there is linearity
Mean square model, msq(model), has expected value $\sigma_{Y X}^2$	Mean square model, msq(model), has expected value $\sigma_{Y X}^2 + \beta_1^2 \sum_{i=1}^n (X_i - \bar{X})^2$
Mean square residual, msq(residual), has expected value $\sigma_{Y X}^2$	Mean square residual, msq(residual), has expected value $\sigma_{Y X}^2$
$F = \text{msq(model)}/\text{msq(residual)}$ tends to be close to 1	$F = \text{msq(model)}/\text{msq(residual)}$ tends to be LARGER than 1

Art of Stat: Dependent = Z, Predictor = X

Regression Options:

☐ Find Predicted Value

☐ Show Residuals on Plot

☐ Show Standard Errors & P-values

☐ Confidence Interval for Slope

☐ Confidence/Prediction Interval

☒ ANOVA Table

Model Summary:

Statistic	Value
Correlation Coefficient r	0.999
Coefficient of Determination r^2	99.8%
Residual Standard Deviation	0.02807

ANOVA Table:

Source	Df	Sum Sq	Mean Sq	F value	P-value
Regression	1	4.2211	4.22106	5,355.6	<0.0001
Residuals	9	0.0071	0.00079		
Total	10	4.2282			

Residual Standard Deviation = Sqrt of MSQ(residual) = $\sqrt{0.00079} = .02807$
 Coefficient of Determination, $r^2 = R^2 = \text{SSQ(model)}/\text{SSQ(TOTAL)} = 4.2211 / 4.2282 = .998$, or 99.8%
 F-value = Mean Sq(Regression)/Mean Sq(Residuals) = $4.22106/0.00079 = 5335.6$

Nature _____ Population/ Sample _____ Observation/ Data _____ Relationships/ Modeling _____ Analysis/ Synthesis

R. Annotations highlighted in yellow.

```
ANOVA Table: Dependent=z_logwt Predictor=x_age
fit2 <- lm(z_logwt ~ x_age, data=dataset)
anova(fit2)

## Analysis of Variance Table
##
## Response: z_logwt
##           Df Sum Sq Mean Sq F value    Pr(>F)
## x_age      1 4.2211   4.2211   5355.6 0.00000000000008399
## Residuals  9 0.0071   0.0008

summary(fit2)
##
## Call:
## lm(formula = z_logwt ~ x_age, data = dataset)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.04854 -0.01787  0.00400  0.02168  0.03402
##
## Coefficients:
##              Estimate Std. Error t value    Pr(>|t|)
## (Intercept) -2.689255    0.030637  -87.78 0.000000000000164
## x_age        0.195891    0.002677   73.18 0.000000000000840
##
## Residual standard error: 0.02807 on 9 degrees of freedom Residual standard error = Sqrt of MSQ(residual) = √0.0008
## Multiple R-squared:  0.9983, Adjusted R-squared:  0.9981 Multiple R-squared = R² = SSQ(model)/SSQ(TOTAL)
## F-statistic: 5356 on 1 and 9 DF, p-value: 0.0000000000008399 F-statistic = F = MSQ(model)/MSQ(residual)
##                                     = 4.2211/0.0008
```

This output corresponds to the following.

Note – In this example our dependent variable is actually Z, not Y.

Source	df	Sum of Squares	Mean Square
Regression <i>due model</i>	1	$SSR = \sum_{i=1}^n (\hat{Z}_i - \bar{Z})^2 = 4.22063$	$msq(model) = SSR/1$ $= 4.22063/1$ $= 4.22063$ <i>You might see</i> $msq(model) = msr$
Residual <i>due error</i>	$(n-2) = 9$	$SSE = \sum_{i=1}^n (Z_i - \hat{Z}_i)^2 = 0.00705$	$msq(residual) = SSE/(n-2)$ $= 0.00705/9$ $= 0.00078$ <i>You might see</i> $msq(residual) = mse$
Total, corrected	$(n-1) = 10$	$SST = \sum_{i=1}^n (Z_i - \bar{Z})^2 = 4.22768$	

Nature _____ Population/ Sample _____ Observation/ Data _____ Relationships/ Modeling _____ Analysis/ Synthesis

Other information in this output:

- ◆ **R-SQUARED** = [(Sum of squares regression)/(Sum of squares total)]
 = proportion of the “total” that we have been able to explain with the fit
 = “percent of variance explained by the model”
 - ***Be careful!*** As predictors are added to the model, R-SQUARED can only increase. Eventually, we need to “adjust” this measure to take this into account. See ADJUSTED R-SQUARED.
- ◆ We also get an overall F test of the null hypothesis that the simple linear model does not explain significantly more variability in LOGWT than the average LOGWT. $F = \text{MSQ (Regression)} / \text{MSQ (Residual)}$

$$= 4.22063 / 0.0007838$$

$$= 5384.94 \text{ with } df = 1, 9$$

p-value = achieved significance < 0.0001. This is a highly unlikely outcome! → Reject H_0 .
 Conclude that the fitted line explains statistically significantly more of the variability in $Z = \text{LOGWT}$ than is explained by the intercept-only null hypothesis model.

7. Hypothesis Testing

Straight Line Model: $Y = \beta_0 + \beta_1 X$

1) Overall F-Test

Research Question: Does the fitted model, the \hat{Y} , explain significantly more of the total variability of the Y about \bar{Y} than does \bar{Y} ?

A bit of clarification here, in case you're wondering. When the null hypothesis is true, at least two things happen: (1) $\beta_1 = 0$ and (2) the correct model (the null one) says $Y = \beta_0 + \text{error}$. In this situation, the least squares estimate of β_0 turns out to be \bar{Y} (that seems reasonable, right?)

Assumptions: As before.

H_0 and H_A :

$$H_0: \beta_1 = 0$$

$$H_A: \beta_1 \neq 0$$

Test Statistic:

$$F = \frac{msq(\text{regression})}{msq(\text{residual})}$$

$$df = 1, (n - 2)$$

Evaluation rule:

When the null hypothesis is true, the value of F should be close to 1. Alternatively, when $\beta_1 \neq 0$, the value of F will be LARGER than 1.

Thus, our p-value calculation answers: "What are the chances of obtaining our value of the F or one that is larger if we believe the null hypothesis that $\beta_1 = 0$?"

Calculations:

For our data, we obtain p-value =

$$\text{pr} \left[F_{1, (n-2)} \geq \mid \frac{msq(\text{model})}{msq(\text{residual})} \mid b_1=0 \right] = \text{pr} [F_{1,9} \geq 5384.94] < .0001$$

Evaluate:

Assumption of the null hypothesis that $\beta_1 = 0$ has led to an extremely unlikely outcome (F-statistic value of 5394.94), with chances of being observed less than 1 chance in 10,000. The null hypothesis is rejected.

Interpret:

We have learned that, at least, the fitted straight line model does a much better job of explaining the variability in $Z = \text{LOGWT}$ than a model that allows only for the average LOGWT.

... later ... (BIOSTATS 640, Intermediate Biostatistics), we'll see that the analysis does not stop here ...

Art of Stat: Dependent = Z, Predictor = X

Regression Options:

- ☐ Find Predicted Value
- ☐ Show Residuals on Plot
- ☐ Show Standard Errors & P-values
- ☐ Confidence Interval for Slope
- ☐ Confidence/Prediction Interval
- ☒ ANOVA Table

Model Summary:

Statistic	Value
Correlation Coefficient r	0.999
Coefficient of Determination r^2	99.8%
Residual Standard Deviation	0.02807

ANOVA Table:

Source	Df	Sum Sq	Mean Sq	F value	P-value
Regression	1	4.2211	4.22106	5,355.6	<0.0001
Residuals	9	0.0071	0.00079		
Total	10	4.2282			

$$F\text{-value} = \text{Mean Sq(Regression)} / \text{Mean Sq(Residuals)} = 4.22106 / 0.00079 = 5335.6$$

R

```
# ANOVA Table: Dependent=z_logwt Predictor=x_age
fit2 <- lm(z_logwt ~ x_age, data=dataset)
anova(fit2)

## Analysis of Variance Table
##
## Response: z_logwt
##      Df Sum Sq Mean Sq F value    Pr(>F)
## x_age   1  4.2211   4.2211  5355.6 0.00000000000008399
## Residuals  9  0.0071   0.0008

summary(fit2)
--- some output not shown ---

## F-statistic: 5356 on 1 and 9 DF, p-value: 0.00000000000008399
```

$$F = \text{MSQ(model)} / \text{MSQ(residual)}$$

$$= [4.2211] / [0.0008]$$



2) Test of the Slope, β_1

Notes -

The overall F test and the test of the slope are equivalent. The test of the slope uses a t-score approach to hypothesis testing. It can be shown that $\{ \text{t-score for slope} \}^2 = \{ \text{overall F} \}$

Research Question: Is the slope $\beta_1 = 0$?

Assumptions: As before.

H_0 and H_A :

$$H_0 : \beta_1 = 0$$

$$H_A : \beta_1 \neq 0$$

Test Statistic:

To compute the t-score, we need an estimate of the standard error of $\hat{\beta}_1$

$$SE(\hat{\beta}_1) = \sqrt{msq(residual) \left[\frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]}$$

Recall what we mean by a t-score:

$t=73.18$ says “the estimated slope is estimated to be 73.18 standard error units away from the null hypothesis expected value of zero”.

Check that $\{t\text{-score}\}^2 = \{\text{Overall } F\}$:

$[73.18]^2 = 5355.3124$ which is close.

Evaluation rule:

When the null hypothesis is true, the value of t should be close to zero.
Alternatively, when $\beta_1 \neq 0$, the value of t will be DIFFERENT from 0.

Here, our p-value calculation answers: “Under the assumption of the null hypothesis that $\beta_1 = 0$, what were our chances of obtaining a t-statistic value 73.18 standard error units away from its null hypothesis expected value of zero”?

Calculations:

For our data, we obtain p-value =

$$2pr\left[t_{(n-2)} \geq \left| \frac{\hat{\beta}_1 - 0}{se(\hat{\beta}_1)} \right| \right] = 2pr[t_9 \geq 73.38] < .0001$$

Evaluate:

Under the null hypothesis that $\beta_1 = 0$, the chances of obtaining a t-score value that is 73.18 or more standard error units away from the expected value of 0 is less than 1 chance in 10,000.

Interpret:

The inference is the same as that for the overall F test. The fitted straight line model does a statistically significantly better job of explaining the variability in LOGWT than the sample mean.

3) Test of the Intercept, β_0

This addresses the question: Does the straight-line relationship passes through the origin? It is rarely of interest.

Research Question: Is the intercept $\beta_0 = 0$?

Assumptions: As before.

H_0 and H_A :

$$H_0 : \beta_0 = 0$$

$$H_A : \beta_0 \neq 0$$

Test Statistic:

To compute the t-score for the intercept, we need an estimate of the standard error of $\hat{\beta}_0$

$$SE(\hat{\beta}_0) = \sqrt{msq(residual) \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]}$$

Our t-score is therefore:

$$t - score = \left[\frac{(observed) - (expected)}{se(expected)} \right] = \left[\frac{(\hat{\beta}_0) - (0)}{se(\hat{\beta}_0)} \right]$$

$$df = (n - 2)$$

Art of Stat: Dependent = Z, Predictor = X

Regression Options:

☐ Find Predicted Value

☐ Show Residuals on Plot

☒ Show Standard Errors & P-values

☐ Confidence Interval for Slope

☐ Confidence/Prediction Interval

☐ ANOVA Table

Linear Regression Equation:

Parameter	Estimate	Standard Error	t Statistic	P-value
Intercept	-2.689	0.03064	-87.78	<0.0001
Slope (X)	0.1959	0.002677	73.18	<0.0001

$$t \text{ Statistic} = \text{Estimate} / \text{Standard Error} = [-2.689] / [0.03064] = -87.78$$

R

```
# TEST OF INTERCEPT: Dependent=z_logwt Predictor=x_age
fit2 <- lm(z_logwt ~ x_age, data=dataset)
summary(fit2)

-- some output not shown --
##
## Coefficients:
##              Estimate Std. Error t value      Pr(>|t|)
## (Intercept) -2.689255   0.030637  -87.78 0.0000000000000164
## x_age        0.195891   0.002677   73.18 0.0000000000000840
```

Here, $t = -87.78$ says “the estimated intercept is estimated to be 87.78 standard error units away from its null hypothesis expected value of zero”.

Evaluation rule:

When the null hypothesis is true, the value of t should be close to zero.
Alternatively, when $\beta_0 \neq 0$, the value of t will be DIFFERENT from 0.

Our p-value calculation answers: “Under the assumption of the null hypothesis that $\beta_0 = 0$, what were our chances of obtaining a t-statistic value 87.78 standard error units away from its null hypothesis expected value of zero”?



Calculations:

p-value =

$$2pr \left[t_{(n-2)} \geq \left| \frac{\hat{\beta}_0 - 0}{s\hat{e}(\hat{\beta}_0)} \right| \right] = 2pr [t_9 \geq 87.78] < .0001$$

Evaluate:

Under the null hypothesis that the line passes through the origin, that $\beta_0 = 0$, the chances of obtaining a t-score value that is 87.78 or more standard error units away from the expected value of 0 is less than 1 chance in 10,000, again prompting statistical rejection of the null hypothesis.

Interpret:

The inference is that there is statistically significant evidence that the straight-line relationship between $Z=\text{LOGWT}$ and $X=\text{AGE}$ does not pass through the origin.

8. Confidence Interval Estimation

Straight Line Model: $Y = \beta_0 + \beta_1 X$

The confidence intervals here have the usual 3 elements (for review, see again Units 8, 9 & 10):

- 1) Best single guess (estimate)
- 2) Standard error of the best single guess (SE[estimate])
- 3) Confidence coefficient: This will be a percentile from the Student t distribution with $df=(n-2)$

We might want confidence interval estimates of the following 4 parameters:

- (1) Slope
- (2) Intercept
- (3) Mean of subset of population for whom $X=x_0$
- (4) Individual response for person for whom $X=x_0$

1) SLOPE

$$\text{estimate} = \hat{\beta}_1$$

$$s\hat{e}(\hat{b}_1) = \sqrt{\text{msq}(\text{residual}) \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2}} = \sqrt{(\text{mse}) \frac{1}{\sum_{i=1}^n (X_i - \bar{X})^2}}$$

2) INTERCEPT

$$\text{estimate} = \hat{\beta}_0$$

$$s\hat{e}(\hat{b}_0) = \sqrt{\text{msq}(\text{residual}) \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]} = \sqrt{(\text{mse}) \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]}$$

3) MEAN at $X=x_0$ estimate = $\hat{Y}_{X=x_0} = \hat{\beta}_0 + \hat{\beta}_1 x_0$

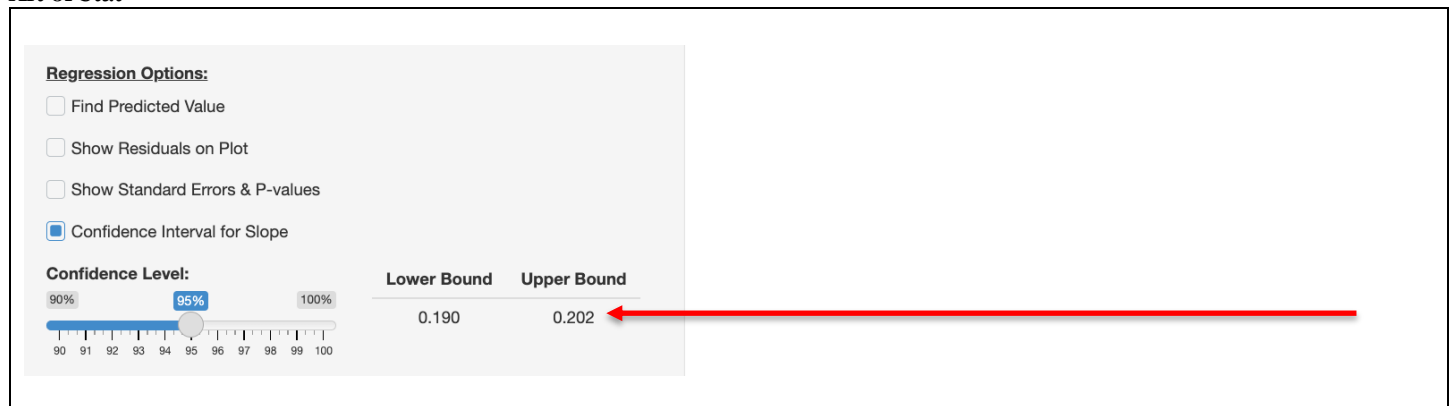
$$\hat{s}_e = \sqrt{\text{msq}(\text{residual}) \left[\frac{1}{n} + \frac{(x_0 - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]} = \sqrt{(\text{mse}) \left[\frac{1}{n} + \frac{(x_0 - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]}$$

4) INDIVIDUAL with $X=x_0$ estimate = $\hat{Y}_{X=x_0} = \hat{\beta}_0 + \hat{\beta}_1 x_0$

$$\hat{s}_e = \sqrt{\text{msq}(\text{residual}) \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]} = \sqrt{(\text{mse}) \left[1 + \frac{1}{n} + \frac{(x_0 - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \right]}$$

1) Confidence Interval for SLOPE Z=LOGWT to X=AGE.

Art of Stat



R

```
fit2 <- lm(z_logwt ~ x_age, data=dataset)
confint(fit2, level=.95)
```

```
##               2.5 %       97.5 %
## (Intercept) -2.7585602 -2.6199489
## x_age       0.1898356  0.2019462
```

With 95% confidence, the slope is estimated to be between 0.19 and 0.20

By Hand

95% Confidence Interval for the Slope, β_1

$$1) \text{ Best single guess (estimate)} = \hat{\beta}_1 = 0.19589$$

$$2) \text{ Standard error of the best single guess (SE[estimate])} = se(\hat{\beta}_1) = 0.00268$$

$$3) \text{ Confidence coefficient} = 97.5^{\text{th}} \text{ percentile of Student } t = t_{.975, df=9} = 2.26$$

95% Confidence Interval for Slope β_1 = Estimate \pm (confidence coefficient) * SE

$$= 0.19589 \pm (2.26)(0.00268)$$

$$= (0.1898, 0.2019)$$

2) Confidence Interval for INTERCEPT

Z=LOGWT to X=AGE.

Art of Stat

--- currently not available ---

R

```
fit2 <- lm(z_logwt ~ x_age, data=dataset)
confint(fit2, level=.95)

##                2.5 %      97.5 %
## (Intercept) -2.7585602 -2.6199489 With 95% confidence, the intercept is estimated to be between -2.76 and -2.62
## x_age       0.1898356  0.2019462
```

By Hand

$$1) \text{ Best single guess (estimate)} = \hat{\beta}_0 = -2.68925$$

$$2) \text{ Standard error of the best single guess (SE[estimate])} = se(\hat{\beta}_0) = 0.03064$$

$$3) \text{ Confidence coefficient} = 97.5^{\text{th}} \text{ percentile of Student } t = t_{.975, df=9} = 2.26$$

95% Confidence Interval for Slope β_0 = Estimate \pm (confidence coefficient) * SE

$$= -2.68925 \pm (2.26)(0.03064)$$

$$= (-2.7585, -2.6200)$$

3) Confidence Interval for MEANS Z=LOGWT to X=AGE.

Art of Stat

NOTE: Art of Stat currently provides confidence intervals for the mean at single values of X ONLY



R

```
# Confidence Intervals for MEAN at each value of X
Dependent=z_logwt Predictor=x_age

# values of x: age=6,7,8,9,10,11,12,13,14,15,16
mydata <- data.frame(x_age = seq(from=6, to=16, by=1))

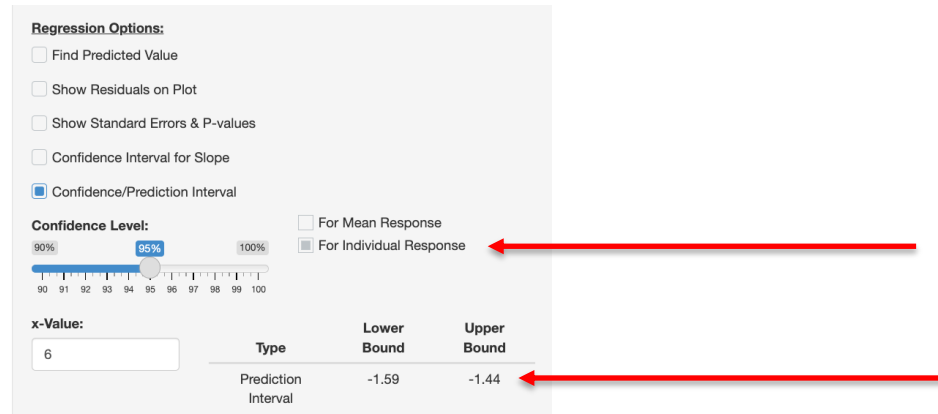
# Confidence Interval Estimates of Means
estimated_mean <- predict(fit2, newdata = mydata, interval = "confidence")
out_mean <- cbind(mydata,estimated_mean) # cbind( ) appends age for readability
out_mean
```

	x_age	fit	lwr	upr
1	6	-1.51390909	-1.54973251	-1.47808568
2	7	-1.31801818	-1.34889408	-1.28714228
3	8	-1.12212727	-1.14852155	-1.09573300
4	9	-0.92623636	-0.94889308	-0.90357965
5	10	-0.73034545	-0.75042849	-0.71026242
6	11	-0.53445455	-0.55360297	-0.51530612
7	12	-0.33856364	-0.35864667	-0.31848060
8	13	-0.14267273	-0.16532944	-0.12001601
9	14	0.05321818	0.02682391	0.07961246
10	15	0.24910909	0.21823319	0.27998499
11	16	0.44500000	0.40917659	0.48082341

4) Confidence Interval for INDIVIDUAL PREDICTIONS Z=LOGWT to X=AGE.

Art of Stat

NOTE: Art of Stat currently provides confidence intervals for the mean at single values of X ONLY



R

Confidence Intervals for INDIVIDUAL PREDICTION at each value of X Dependent=z_logwt Predictor=x_age

```
estimated_individual <- predict(fit2, newdata=mydata, interval = "prediction")
out_individual <- cbind(mydata,estimated_individual)
out_individual
```

	x_age	fit	lwr	upr
1	6	-1.51390909	-1.58682410	-1.44099408
2	7	-1.31801818	-1.38863407	-1.24740230
3	8	-1.12212727	-1.19090183	-1.05335271
4	9	-0.92623636	-0.99366491	-0.85880782
5	10	-0.73034545	-0.79695334	-0.66373757
6	11	-0.53445455	-0.60078662	-0.46812247
7	12	-0.33856364	-0.40517152	-0.27195575
8	13	-0.14267273	-0.21010127	-0.07524418
9	14	0.05321818	0.01555638	0.12199274
10	15	0.24910909	0.17849320	0.31972498
11	16	0.44500000	0.37208499	0.51791501

9. Introduction to Correlation

Definition of Correlation

A correlation coefficient is a measure of the association between two paired random variables (e.g. height and weight).

The **Pearson product moment correlation**, in particular, is a measure of the strength of the *straight-line* relationship between the two random variables.

Another correlation measure (not discussed here) is the **Spearman correlation**. It is a measure of the strength of the *monotone increasing (or decreasing)* relationship between the two random variables. The Spearman correlation is a non-parametric (meaning model free) measure. It is introduced in BIOSTATS 640, *Intermediate Biostatistics*.

Formula for the Pearson Product Moment Correlation ρ

- Population product moment correlation = ρ
- Sample based estimate = r .
- Some preliminaries:

(1) Suppose we are interested in the correlation between X and Y

$$(2) \text{cov}\hat{(X,Y)} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1)} = \frac{S_{xy}}{(n-1)} \quad \text{This is the covariance}(X,Y)$$

$$(3) \text{var}\hat{(X)} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n-1)} = \frac{S_{xx}}{(n-1)} \quad \text{and similarly}$$

$$(4) \text{var}\hat{(Y)} = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{(n-1)} = \frac{S_{yy}}{(n-1)}$$

Formula for Estimate of Pearson Product Moment Correlation from a Sample

$$\hat{\rho} = r = \frac{\text{cov}(\hat{x}, \hat{y})}{\sqrt{\text{var}(\hat{x})\text{var}(\hat{y})}}$$

$$= \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}}$$

If you absolutely have to do it by hand, an equivalent (more calculator/excel friendly formula) is

$$\hat{\rho} = r = \frac{\sum_{i=1}^n x_i y_i - \frac{\left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n}}{\sqrt{\left[\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}\right] \left[\sum_{i=1}^n y_i^2 - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n}\right]}}$$

- The correlation r can take on values **between 0 and 1 only**
- Thus, the correlation coefficient is said to be **dimensionless** – it is independent of the units of x or y .
- **Sign** of the correlation coefficient (positive or negative) = **Sign** of the estimated slope $\hat{\beta}_1$.

There is a relationship between the slope of the straight line, $\hat{\beta}_1$, and the estimated correlation r .

Relationship between slope $\hat{\beta}_1$ and the sample correlation r

Tip! This is very handy...

Because $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$ and $r = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}}$

A little algebra reveals that

$$r = \left[\frac{\sqrt{S_{xx}}}{\sqrt{S_{yy}}} \right] \hat{\beta}_1$$

Thus, beware!!!

- It is possible to have a very large (positive or negative) r might accompanying a very non-zero slope, inasmuch as
 - A very large r might reflect a very large S_{xx} , all other things equal
 - A very large r might reflect a very small S_{yy} , all other things equal.

10. Hypothesis Test of Correlation

The null hypothesis of zero correlation is equivalent to the null hypothesis of zero slope.

Research Question: Is the correlation $\rho = 0$? Is the slope $\beta_1 = 0$?

Assumptions: As before.

H_0 and H_A :

$$H_0 : \rho = 0$$

$$H_A : \rho \neq 0$$

Test Statistic:

A little algebra (not shown) yields a very nice formula for the t-score that we need.

$$t\text{-score} = \left[\frac{r \sqrt{(n-2)}}{\sqrt{1-r^2}} \right]$$

$$df = (n-2)$$

We can find this information in our output. Recall the first example and the model of Z=LOGWT to X=AGE:

Art of Stat: The Pearson Correlation, r , is the square root of “Coefficient of Determination, r^2 ”

Model Summary:

Statistic	Value
Correlation Coefficient r	0.999
Coefficient of Determination r^2	99.8%
Residual Standard Deviation	0.02807

R: The Pearson Correlation, r , is the square root of “Multiple R-Squared”

```
fit2 <- lm(z_logwt ~ x_age, data=dataset)
summary(fit2)

-- some output not shown --
Residual standard error: 0.02807 on 9 degrees of freedom
Multiple R-squared: 0.9983, Adjusted R-squared: 0.9981 Pearson Correlation,  $r = \sqrt{0.9983} = 0.9991$ 
F-statistic: 5356 on 1 and 9 DF, p-value: 0.0000000000008399
```

Alternatively, you can get the Pearson Correlation, r using the following:
`sqrt(summary(fit2)$r.squared)`

```
[1] 0.9991608
```

Terrific! We can also extract the value of the Pearson Correlation, r , from the T-statistic for Null: Slope = 0
 Substitution into the formula for the t-score yields

$$t\text{-score} = \left[\frac{r \sqrt{(n-2)}}{\sqrt{1-r^2}} \right] = \left[\frac{.9991\sqrt{9}}{\sqrt{1-.9983}} \right] = \left[\frac{2.9974}{.0412} \right] = 72.69$$

Note: The value .9991 in the numerator is $r = \sqrt{R^2} = \sqrt{.9983} = .9991$

This is very close to the value of the t-score that was obtained for testing the null hypothesis of zero slope. The discrepancy is probably rounding error. I did the calculations on my calculator using 4 significant digits.

Nature _____ Population/
Sample _____ Observation/
Data _____ Relationships/
Modeling _____ Analysis/
Synthesis