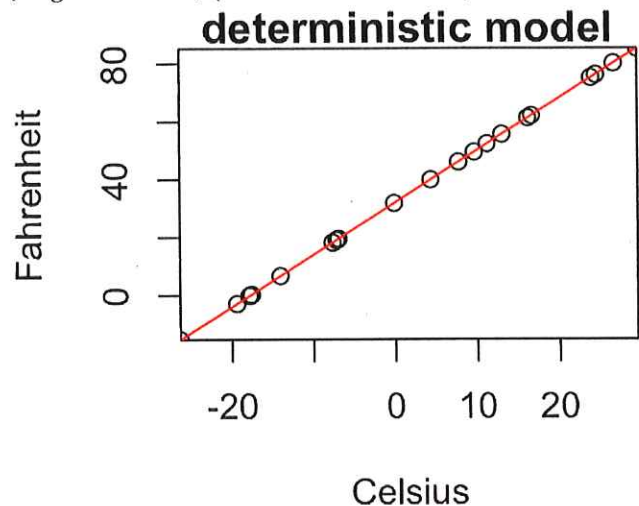Simple Linear Regression: Modelling Linear Relationship

After discovering covariance between bivariate data, you will likely want to know how to describe/express the relationship. Successful descriptions can then be used to model expected results on one of the variables deemed the **response variable, Y**, based only off the **predictor variable, X**.

The steps are:
1. Collect bivariate $X$ and $Y$ variables from historic events. This data set will serve as "training" to understand the linear relationship that exists between the variables.
2. Develop a mathematical expression/equation to transform an particular/hypothetical $X$ into an **expected/estimated** $Y$.

Consider, "bivariate" data expressing temperature in $°C$, *degrees Celsius*, (call this the $X$ variable) and then in $°F$, *degrees Fahrenheit*, (call this the $Y$ variable).

| | Celsius | Fahrenheit |
|---|---|---|
| 1 | -14.0694802 | 6.6749356 |
| 2 | -7.6725660 | 18.1893812 |
| 3 | 4.3712018 | 39.8681632 |
| 4 | 24.4924674 | 76.0864413 |
| 5 | -17.8990841 | -0.2183514 |
| 6 | 23.9033811 | 75.0260860 |
| 7 | 26.6805161 | 80.0249290 |
| 8 | 9.6478675 | 49.3661616 |
| 9 | 7.7468426 | 45.9443167 |
| 10 | -26.2928238 | -15.3270828 |
| 11 | -17.6415255 | 0.2452541 |
| 12 | -19.4065948 | -2.9318707 |



deterministic model — scatterplot of Fahrenheit vs Celsius

What do you notice about the scatterplot?
Do you see how all the points fall exactly on the line? This is called a deterministic model; since all points fall exactly on the line we could perfectly predict where no points exist.

The linear equation will follow the **deterministic model's** form:

$$Y_i = \beta_0 + \beta_1 X_i, \quad i = 1, 2, \dots, n$$

where $\beta_0$ is the y-intercept (the $Y$ value when $X = 0$) and $\beta_1$ is the slope (rate of change in $Y$ with respect to $X$).

In a deterministic model only two sample points are all that is required to find the model:

$$\beta_1 = \frac{rise}{run} = \frac{y_2 - y_1}{x_2 - x_1}, \text{ then } \beta_0 = y_1 - \beta_1 x_1$$

Of course in the "real" world we often lack the ability to measure variables with deterministic precision. We expect response and or measurement bias in our observations. Additionally, when dealing with random variables we know that our observations may lack consistency not due to any bias at all!

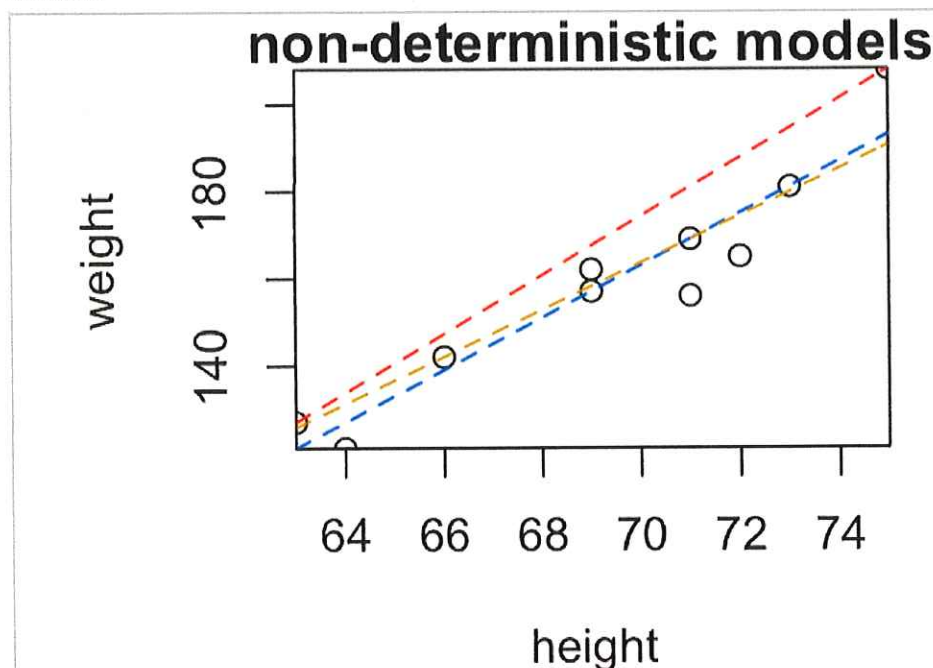$\beta_1 \cong \dfrac{6.675 - 18.189}{-14.069 - -7.673} \cong 1.800$   $\beta_0 \cong 6.675 - (1.8)(-14.069) \cong 32.$

$(°F) = 32 + 1.8(°C)$   $\Rightarrow$ for each Additional $1°C$, $°F$ goes up $1.8$ degrees $\leftarrow$ Slope, $\beta_1$

Stat323©ScottRobison2017

$\Rightarrow$ At $0°C$ $\Rightarrow$ $32°F$ $\leftarrow$ $\beta_0$, y intercept.

Let's return to our example of ten student's height and weight, and try to select two points from the data set then create linear models…

## non-deterministic models



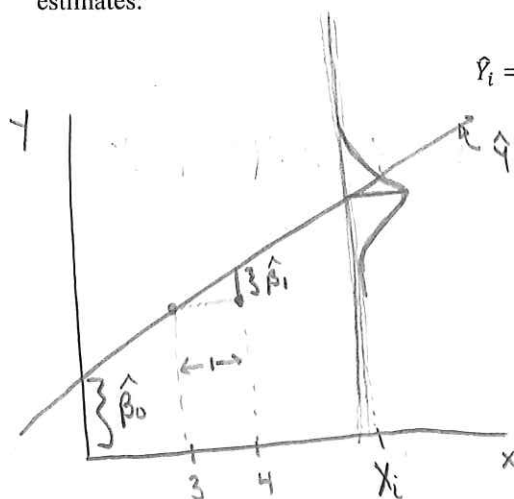weight (y-axis: 140, 180)

height (x-axis: 64 66 68 70 72 74)

So which model of a non-deterministic data set is best?...

The **probabilistic model** appears to the same as the deterministic model, however, it includes an addition $\varepsilon_i$ term:

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \qquad i = 1, 2, \dots, n, \qquad \text{where } \varepsilon_i \sim Norm(\beta_0 + \beta_1 X_i, \sigma^2)$$

The probabilistic model can also be written this way, to "hide" the $\varepsilon_i$ term by admitting the following values are estimates:

$$\hat{Y}_i = \widehat{\beta_0} + \widehat{\beta_1} X_i, \qquad i = 1, 2, \dots, n$$



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

$$Y_i \sim Norm\left(\hat{\beta_0} + \hat{\beta_1} X_i, \sigma^2\right)$$

($Y$ is random variable)

OR...

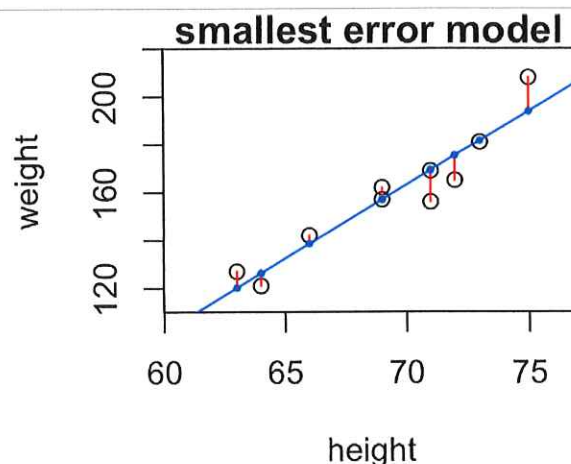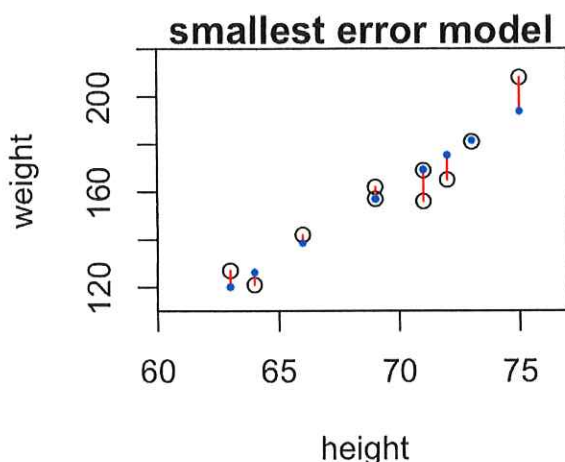$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

$$\varepsilon_i \sim Norm(0, \sigma^2)$$

($Y$ is a conditional constant, that is measured inaccurately)

$\hat{Y} = \widehat{\beta_0} + \widehat{\beta_1} X$; The Least-Squares Estimate Model

Let $\hat{Y}$ be the probabilistic model that is the **closest** "overall" to each sample point, meaning the set of sample points $(X_i, Y_i)'s$ that are respectively **closest** to the $(X_i, \hat{Y}_i)'s$. Then we define the difference between $Y_i$ and $\hat{Y}_i$ to be $\varepsilon_i$.

Then $\varepsilon_i = Y_i - \hat{Y}_i$ (residuals/errors/residual errors), we are setting $\sum_{i=1}^n \varepsilon_i = \sum_{i=1}^n (Y_i - \hat{Y}_i) = 0$ so we can find the model with the **least** overall error!



One complication that comes up due to setting the sum of errors equal to zero is that you have now clearly made the signs of some of the errors negative.

To overcome this we square the individual error terms and then discuss the SUM of SQUARED ERRORS,
$$SSE = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - (\widehat{\beta_0} + \widehat{\beta_1} X_i))^2 \qquad Y_i \to observed \ Value \qquad \hat{Y}_i \to predicted \ Value$$

We wish to estimate the probabilistic model in such a way that the square of these vertical distances is as small as possible, a method that invokes **least-squares estimation**. Consider the sum of the squared distances/errors, SSE, each bivariate data point lies away from the imaginary linear line, $\hat{Y} = \widehat{\beta_0} + \widehat{\beta_1} X$.

We need to **minimize** SSE with respect to $\widehat{\beta_0}$ then with respect to $\widehat{\beta_1}$ by finding then setting the partial derivatives equal to zero. $\frac{\delta SSE}{\delta \widehat{\beta_i}}$, where $i = 0, 1$

We will see:

The least-squares estimate of the Y-intercept of the model is:

$$\widehat{\beta_0} = \bar{Y} - \widehat{\beta_1}\bar{X}$$

The least-squares estimate of the slope of the model is:

$$\widehat{\beta_1} = \frac{S_{XY}}{S_{XX}} = \frac{S_{XY}}{S_X S_X} = \frac{S_{XY}}{S_X^2} = r\frac{S_y}{S_X} = \frac{\sum_{i=1}^n [(X_i - \bar{X})(Y_i - \bar{Y})]}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i Y_i) - n\bar{X}\bar{Y}}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i Y_i) - n\bar{X}\bar{Y}}{\sum_{i=1}^n X_i^2 - n\bar{X}^2}$$

Stat323©ScottRobison2017

$$\frac{\delta SSE}{\delta \hat{\beta_0}} = \frac{\partial}{\partial \hat{\beta_0}}\left[\sum(y_i - (\hat{\beta_0} + \hat{\beta_1}x_i))^2\right]$$

need to expand $\left[y_i - (\hat{\beta_0} + \hat{\beta_1}x_i)\right]^2$ term.

expansion:

$(y_i - \hat{\beta_0} - \hat{\beta_1}x_i)(y_i - \hat{\beta_0} - \hat{\beta_1}x_i)$

However I will be taking the derivative w.r.t $\hat{\beta_0}$

So... I only care about terms with $\hat{\beta_0}$ in them...

Only 5 terms matter...

$\propto \frac{\partial}{\partial \hat{\beta_0}}\sum\left[-y_i\hat{\beta_0} + \hat{\beta_0}^2 + \hat{\beta_0}\hat{\beta_1}x_i - y_i\hat{\beta_0} + \hat{\beta_0}\hat{\beta_1}x_i\right] = \frac{\partial}{\partial \hat{\beta_0}}\sum\left[-2y_i\hat{\beta_0} + \hat{\beta_0}^2 + 2\hat{\beta_0}\hat{\beta_1}x_i\right]$

Sum of ders = der of sums

$$= \sum\left[-2y_i + 2\hat{\beta_0} + 2\hat{\beta_1}x_i\right] = -2\sum y_i + 2n\hat{\beta_0} + 2\hat{\beta_1}\sum x_i$$

find min/max   set = 0! Solve for $\hat{\beta_0}$!

$$-2\sum y_i + 2n\hat{\beta_0} + 2\hat{\beta_1}\sum x_i = 0 \implies \frac{2n\hat{\beta_0}}{2n} = \frac{2\sum y_i}{2n} - \frac{2\hat{\beta_1}\sum x_i}{2n}$$

$$\implies \hat{\beta_0} = \frac{\sum y_i}{n} - \hat{\beta_1}\frac{\sum x_i}{n} = \boxed{\overline{Y} - \hat{\beta_1}\overline{X} = \hat{\beta_0}}$$

ensure it is a **min** (want to min Error to have least-square estimate!)

$$\frac{\partial^2(SSE)}{\partial \hat{\beta_0}^2} \propto \frac{\partial}{\partial \hat{\beta_0}}\left(-2\sum y_i + 2n\hat{\beta_0} + \hat{\beta_1}\sum x_i\right) = 2n > 0 \text{ since } n \geq 1$$

Positive!

So by $2^{nd}$ derivative test. $\hat{\beta_0} = \overline{Y} - \hat{\beta_1}\overline{X}$ minimizes $\underline{SSE}$!

$$\frac{\delta SSE}{\delta \hat{\beta}_1} = \frac{\partial}{\partial \hat{\beta}_1}\left[\sum(y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2\right]$$

Need to expand $\left[y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)\right]^2$ term ...

Consider expansion:

$$(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)$$

only five terms will matter!

and only three like terms

However, I know I will Be taking the derivative wrt. $\hat{\beta}_1$. So ... I only care about terms that will have $\hat{\beta}_1$ in them. ignore the rest? ...

$$\frac{\delta SSE}{\delta \hat{\beta}_1} \propto \frac{\partial}{\partial \hat{\beta}_1}\left[\sum(\hat{\beta}_1^2 x_i^2 + 2\hat{\beta}_1\hat{\beta}_0 x_i - 2\hat{\beta}_1 x_i y_i)\right]$$

Sum of der is same as der of sum

$$= \sum\left[2\hat{\beta}_1 x_i^2 + 2\hat{\beta}_0 x_i - 2x_i y_i\right] = -2 \cdot \sum(x_i(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i))$$

Set $\frac{\partial SSE}{\partial \hat{\beta}_1} = 0$ to find min or max

$$-2\sum(x_i(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)) = 0 \implies \sum(x_i(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)) = 0$$

$$\implies \sum(x_i y_i - \hat{\beta}_0(x_i) - \hat{\beta}_1 x_i^2) = \sum x_i y_i - \hat{\beta}_0 \sum x_i - \hat{\beta}_1 \sum x_i^2 = 0$$

$$\implies \hat{\beta}_1 \sum x_i^2 = \sum x_i y_i - \hat{\beta}_0 \sum x_i$$

recall $\left(\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}\right)$ Sub in!

$$\implies \hat{\beta}_1 \sum x_i^2 = \sum x_i y_i - (\bar{y} - \hat{\beta}_1 \bar{x})\sum x_i \left(\frac{n}{n}\right)$$

$$\implies \hat{\beta}_1 \sum x_i^2 = \sum x_i y_i - \bar{x}\bar{y}n - \hat{\beta}_1 \bar{x}^2 n \implies \hat{\beta}_1\left(\sum x_i^2 - n\bar{x}^2\right) = \sum x_i y_i - n\bar{x}\bar{y}$$

$$\implies \boxed{\hat{\beta}_1 = \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2}}$$

ensure $\hat{\beta}_1$ minimizes SSE!

$$\frac{d^2 SSE}{d\hat{\beta}_1^2} \propto \frac{d}{d\hat{\beta}_1} \left( \sum \left( 2\hat{\beta}_1 x_i^2 + 2\hat{\beta}_0 x_i - 2 x_i y_i \right) \right)$$

$$\propto \sum \left( 2 x_i^2 \right) = 2 \sum x_i^2 > 0 \quad \text{clearly positive}$$

Sum of squared values doubled....

then by 2nd der. test $\hat{\beta}_1 = \dfrac{\sum x_i y_i - n \bar{X}\bar{Y}}{\sum x_i^2 - n \bar{x}^2} = \dfrac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$

minimizes SSE!

Example 1: Let's look back at our student height and weight data. Found in data file in D2L, using height as the predictor and weight as the response.

$$Weight_i = \beta_0 + \beta_1 height_i + e_i \iff \widehat{weight_i} = \hat{\beta}_0 + \hat{\beta}_1 height_i$$

$height = c(63, 64, 66, 69, 69, 71, 71, 72, 73, 75)$
$weight = c(127, 121, 142, 157, 162, 156, 169, 165, 181, 208)$

(i)     Find the least-squares estimate of the model.

$$\hat{\beta}_1 = \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2} = \frac{(63 \cdot 127) + (64 \cdot 121) + \cdots + (75 \cdot 208) - 10\left[\frac{63+64+\cdots+75}{10}\right]\left[\frac{127+121+\cdots+208}{10}\right]}{(63^2 + 64^2 + \cdots + 75^2) - 10\left[\frac{63+64+\cdots+75}{10}\right]^2}$$

$\hat{\beta}_1 \cong 6.137581$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \cong \left(\frac{127+121+\cdots+208}{10}\right) - 6.1376\left(\frac{63+64+\cdots+75}{10}\right)$$

$\hat{\beta}_0 \cong -266.534395$

R command:

$fit = lm(Y \sim X)$
$fit \$ coef$

$\Rightarrow \hat{y}_i = -266.534 + 6.138 X_i$

$\Rightarrow \widehat{weight}_i = -266.534 + 6.138(height_i)$

(ii)     Interpret the meaning of $\widehat{\beta}_1$ found in part (i).

$\hat{\beta}_1$ for every additional unit increase in X Expect $\hat{\beta}_1$, unit increases in Y.

Here in Context

for every additional inch in height expect a person to weight and additional ~6.14 lbs...

Causation? Does weight cause height? Does height "cause" weight?

clearly Not The only factor involved...
        Confounding / Additional variables!

(iii) Using the estimate of the model in (i), predict the weight a student with a height of 2.5 feet (recall there are 12 inches in a foot). Interpret the meaning of this predicted value. Try again with a height of 67 inches.

? $\hat{y}_i = -266.534 + 6.138(2.5)$ ⊘ ⟹ $2.5(12) = 30$ inches.

⟹ $\hat{y}_i = -266.534 + 6.138(30) \cong -82.394$.

⟹ We expect people who are 30 inches tall to weigh $-82.4$ lbs...

why is this crazy? ⟹ Scope of sampling from X range from $63-75$

No when near value of 30 we don't know the relationship at values like 30 but clearly it is Not the same relationship as it is around $63-75$ inch in height. (extrapolation!)

⟹ $\hat{y}_i = -266.534 + 6.138(67) \cong 144.713$.

we expect people who are 67 inches tall to weigh about 144.7 lbs!

(iv) Find the value of the residual corresponding to the fourth data point: $(X_4 = 69, Y_4 = \cancel{162}\ 157)$

$\hat{y}_i = -266.534 + 6.138(69) \cong 156.988$

$y_i = 157$.

$e_i = y_i - \hat{y}_i \cong 157 - 156.988 \cong 0.012$.

* See R-code

Example 2: How strong is the linear relationship between the age of a driver and the distance the driver can see? A research firm (Last Resource, Inc., Bellefonte, PA) collected data on a sample of n = 30 drivers. What can you say about the relationship? What is likely to be X? → Easy to know ?

Y? → Interest!

$$\widehat{Distance}_i = 576.6819 - 3.0068 \ Age_i$$

for each year old you can see 3.0068 (feet? / meters? / yard? ... ) <u>less</u>

if you were just born you should be able to see 576.6819 (?)

$Cov(Distance, Age) \cong -1425.862$

$Cor(Distance, Age) \cong -0.8012447.$