

# Correlations and Regression Coefficients

Eli Gurarie

StatR 101 - Lecture 4  
October 22, 2012

October 22, 2012

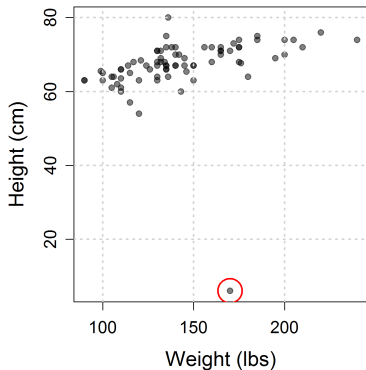
PROFESSIONAL & CONTINUING EDUCATION

UNIVERSITY *of* WASHINGTON



# Scatterplots help identify outliers

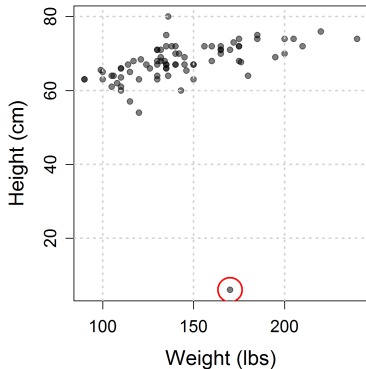
## Data-entry error



Weight: 6 kg?

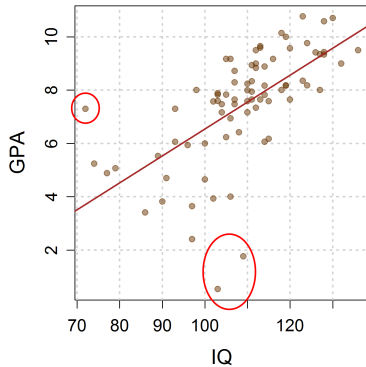
# Scatterplots help identify outliers

## Data-entry error



Weight: 6 kg?

## Informative outliers



Over- and underachievement

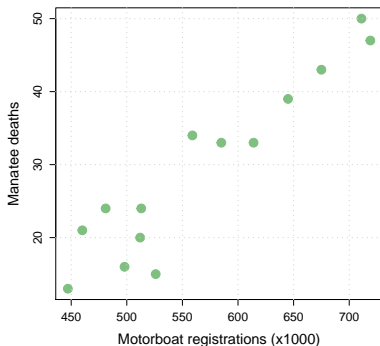
## Correlation ( $r$ )

Is a measure of the **strength** and **direction** of a **linear** relationship

# Manatees and motorboats: Scatterplots

allow us to visually characterize the relationships between *continuous/quantitative* variables.

**Motorboats vs. Manatees**



Identify:

- **direction** (positive/negative),
- **form** (linear/non-linear),
- **strength** (strong/weak)

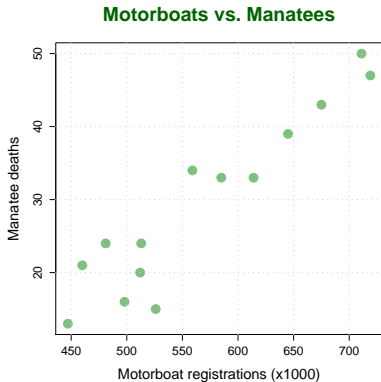
Of a relationship

**R code: scatterplot**

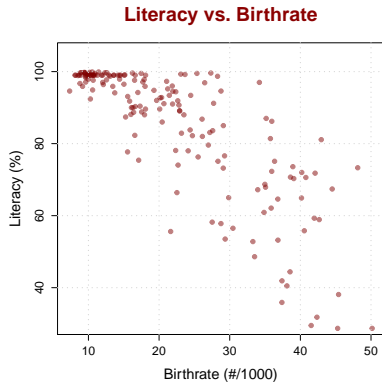
```
plot(Motorboats, Deaths)
```

# Direction of relationship

## Positive relationship

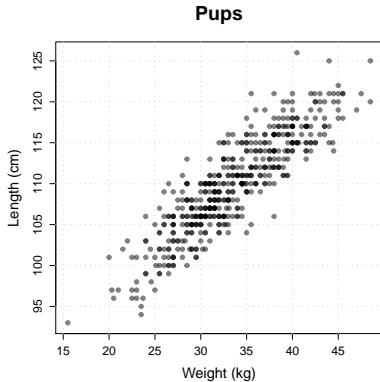


## Negative relationship

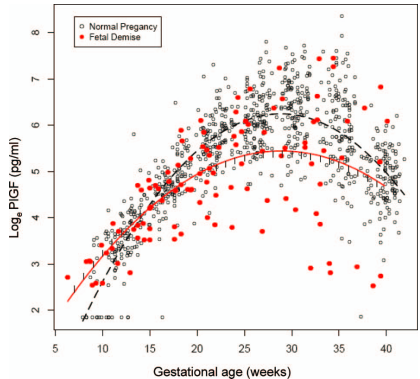


# Form of relationship

## Linear relationship

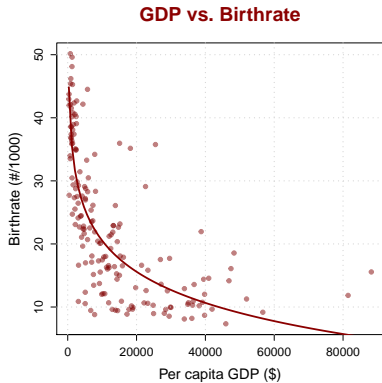


## Non-linear relationship



# Form of relationship

## Non-linear relationship



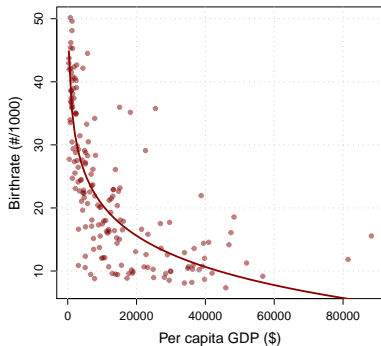


# Form of relationship

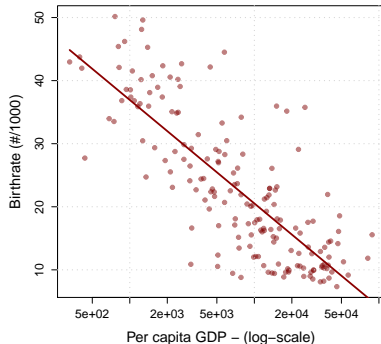
Non-linear relationship

Non-linear ... linearized!

GDP vs. Birthrate



GDP vs. Birthrate

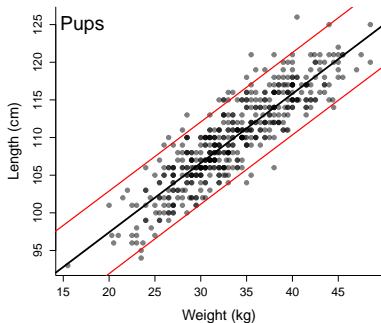


R code: log transformation

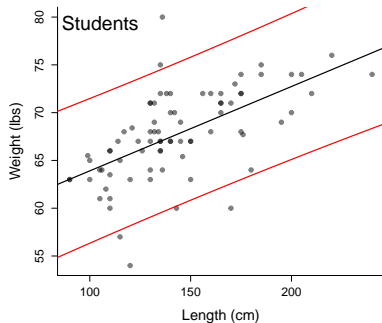
```
plot(GDP, Birthrate, log="x")
```

# Strength of relationship

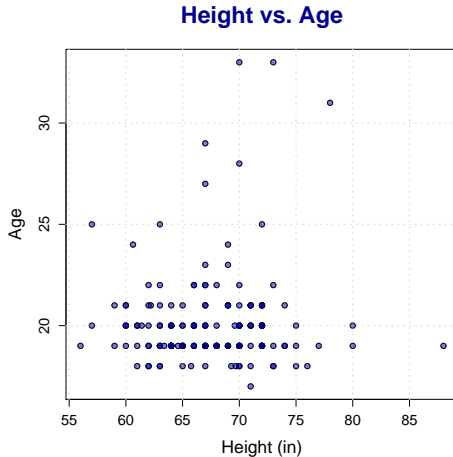
**Stronger relationship**



**Weaker relationship**



# No relationship



Knowing your **height** tells me basically nothing about your **age**.

## Correlation (r)

For any paired sequence of observations:  
 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ :

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

What are the units of the correlation?

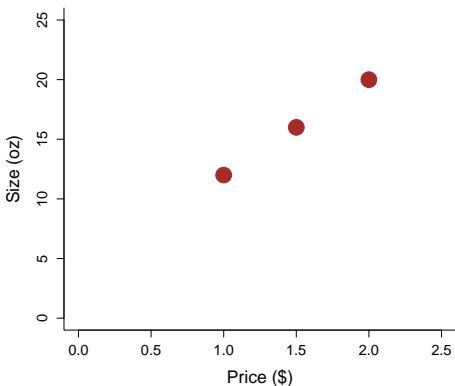
## Correlation: Coffee example

|                 | Price (\$) | Size (oz) |
|-----------------|------------|-----------|
| Tall (small)    | 1.00       | 12        |
| Grande (medium) | 1.50       | 16        |
| Vente (large)   | 2.00       | 20        |



## Correlation: Coffee scatterplot

|                 | Price (\$) | Size (oz) |
|-----------------|------------|-----------|
| Tall (small)    | 1.00       | 12        |
| Grande (medium) | 1.50       | 16        |
| Vente (large)   | 2.00       | 20        |



## Correlation: Coffee calculation

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

|      | Price (\$) | Size (oz) |  |  |
|------|------------|-----------|--|--|
|      | $x$        | $y$       | $\frac{x - \bar{x}}{s_x}$                  | $\frac{y - \bar{y}}{s_y}$                  |
|      |            |           | $\left( \frac{x_i - \bar{x}}{s_x} \right)$ | $\left( \frac{y_i - \bar{y}}{s_y} \right)$ |
|      | 1.00       | 12        |  |  |
|      | 1.50       | 16        |  |  |
|      | 2.00       | 20        |  |  |
| mean |            |           | $\Sigma$                                   |  |
| s.d. |            |           | $r$  |  |

## Correlation: Coffee calculation

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

|      | Price (\$) | Size (oz) |  |  |
|------|------------|-----------|--|--|
|      | $x$        | $y$       | $\frac{x - \bar{x}}{s_x}$                | $\frac{y - \bar{y}}{s_y}$                |
|      |            |           | $\left( \frac{x - \bar{x}}{s_x} \right)$ | $\left( \frac{y - \bar{y}}{s_y} \right)$ |
|      | 1.00       | 12        |  |  |
|      | 1.50       | 16        |  |  |
|      | 2.00       | 20        |  |  |
| mean | 1.5        | 16        | $\Sigma$                                 |  |
| s.d. | 0.5        | 4         | $r$                                      |  |



## Correlation: Coffee calculation

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

|      | Price (\$) | Size (oz) |                           |                           |   |
|------|------------|-----------|---------------------------|---------------------------|---|
|      | $x$        | $y$       | $\frac{x - \bar{x}}{s_x}$ | $\frac{y - \bar{y}}{s_y}$ | $\left( \frac{x - \bar{x}}{s_x} \right) \left( \frac{y - \bar{y}}{s_y} \right)$ |
|      | 1.00       | 12        | -1                        | -1                        |   |
|      | 1.50       | 16        | 0                         | 0                         |   |
|      | 2.00       | 20        | 1                         | 1                         |   |
| mean | 1.5        | 16        |                           | $\Sigma$                  |   |
| s.d. | 0.5        | 4         |                           | $r$                       |   |

## Correlation: Coffee calculation

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

|      | Price (\$) | Size (oz) |                             |                             |   |
|------|------------|-----------|-----------------------------|-----------------------------|---|
|      | $x$        | $y$       | $\frac{x - \bar{x}_i}{s_x}$ | $\frac{y_i - \bar{y}}{s_y}$ | $\left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$ |
|      | 1.00       | 12        | -1                          | -1                          | 1   |
|      | 1.50       | 16        | 0                           | 0                           | 0   |
|      | 2.00       | 20        | 1                           | 1                           | 1   |
| mean | 1.5        | 16        |                             | $\Sigma$                    |   |
| s.d. | 0.5        | 4         |                             | $r$                         |   |

## Correlation: Coffee calculation

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

|      | Price (\$) | Size (oz) |                           |                           |   |
|------|------------|-----------|---------------------------|---------------------------|---|
|      | $x$        | $y$       | $\frac{x - \bar{x}}{s_x}$ | $\frac{y - \bar{y}}{s_y}$ | $\left( \frac{x - \bar{x}}{s_x} \right) \left( \frac{y - \bar{y}}{s_y} \right)$ |
|      | 1.00       | 12        | -1                        | -1                        | 1   |
|      | 1.50       | 16        | 0                         | 0                         | 0   |
|      | 2.00       | 20        | 1                         | 1                         | 1   |
| mean | 1.5        | 16        |                           | $\Sigma$                  | 2   |
| s.d. | 0.5        | 4         |                           | $r$                       | 1   |

## Correlation: Coffee calculation

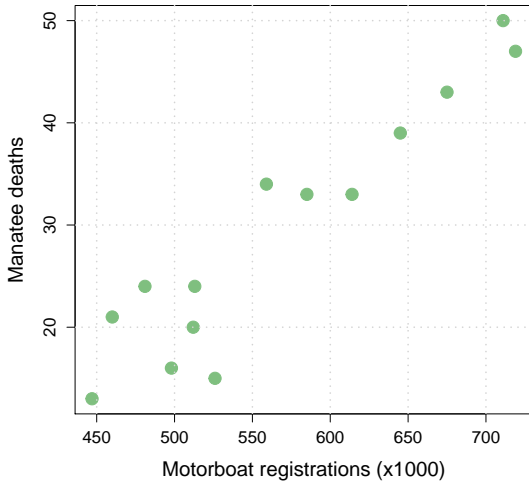
$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

|      | Price (\$) | Size (oz) |                           |                           |   |
|------|------------|-----------|---------------------------|---------------------------|---|
|      | $x$        | $y$       | $\frac{x - \bar{x}}{s_x}$ | $\frac{y - \bar{y}}{s_y}$ | $\left( \frac{x - \bar{x}}{s_x} \right) \left( \frac{y - \bar{y}}{s_y} \right)$ |
|      | 1.00       | 12        | -1                        | -1                        | 1   |
|      | 1.50       | 16        | 0                         | 0                         | 0   |
|      | 2.00       | 20        | 1                         | 1                         | 1   |
| mean | 1.5        | 16        |                           | $\Sigma$                  | 2   |
| s.d. | 0.5        | 4         |                           | $r$                       | 1   |

$r = 1.0$  means PERFECT correlation and a POSITIVE relationship.

## More Correlations

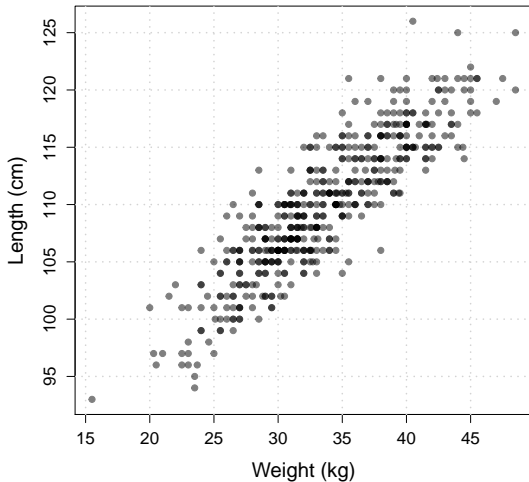
### Motorboats vs. Manatees



$$r = 0.9415$$

## More Correlations

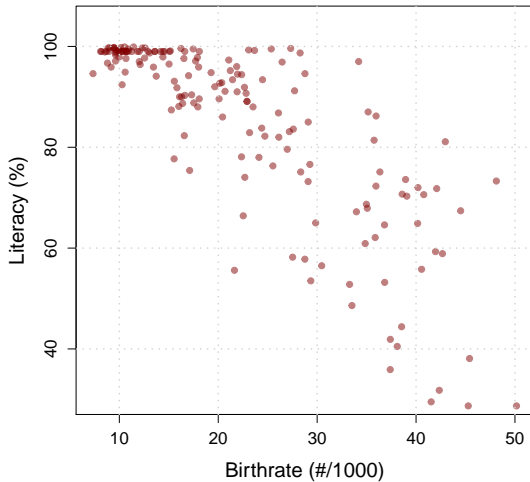
### Pups



$$r = 0.8828$$

## More Correlations

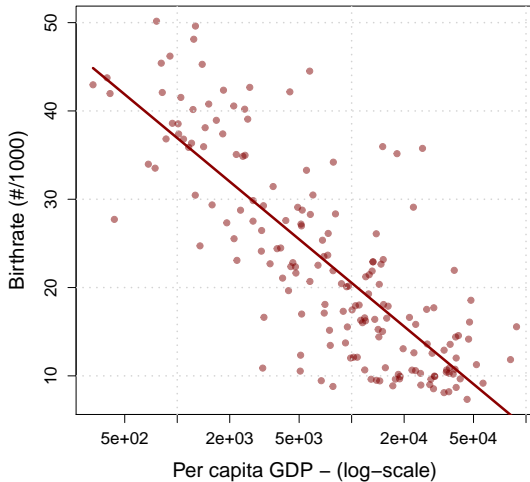
**Literacy vs. Birthrate**



$$r = -0.8138$$

## More Correlations

**GDP vs. Birthrate**

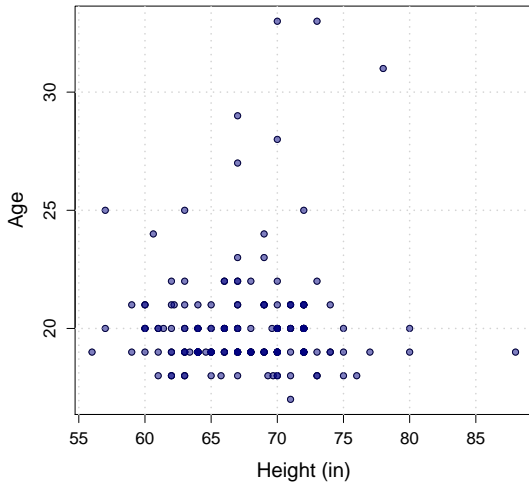


$$r = -0.7986$$



## More Correlations

**Height vs. Age**



$$r = -0.0625$$

## Correlations are...

- unitless and independent of units of measurement;
- between -1 (perfect, negative) and +1 (perfect, positive) with  $r = 0$  meaning no relationship;
- symmetric (no separation between “cause” and “effect”).

## Correlations are...

- unitless and independent of units of measurement;
- between -1 (perfect, negative) and +1 (perfect, positive) with  $r = 0$  meaning no relationship;
- symmetric (no separation between “cause” and “effect”).

## Correlations are...

- unitless and independent of units of measurement;
- between -1 (perfect, negative) and +1 (perfect, positive) with  $r = 0$  meaning no relationship;
- symmetric (no separation between “cause” and “effect”).

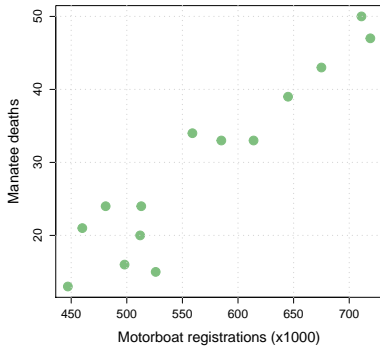
## Correlations are...

- unitless and independent of units of measurement;
- between -1 (perfect, negative) and +1 (perfect, positive) with  $r = 0$  meaning no relationship;
- symmetric (no separation between “cause” and “effect”).

## Why is...

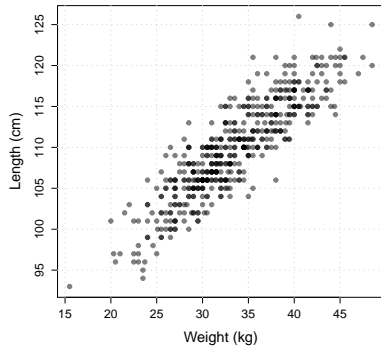
$$r = 0.9415$$

**Motorboats vs. Manatees**



$$r = 0.8828$$

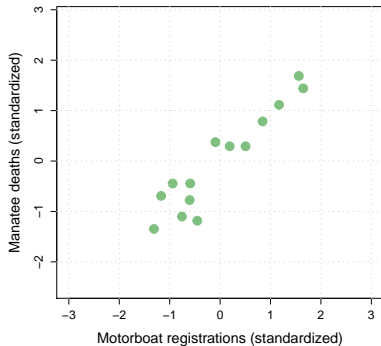
**Pups**



# Standardized scatterplots

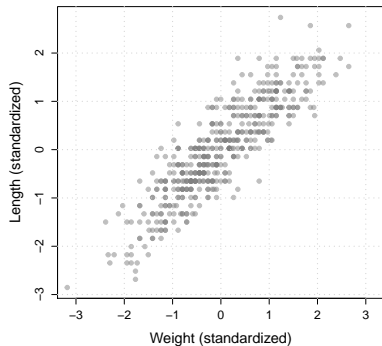
$$r = 0.9415$$

**Motorboats vs. Manatees**

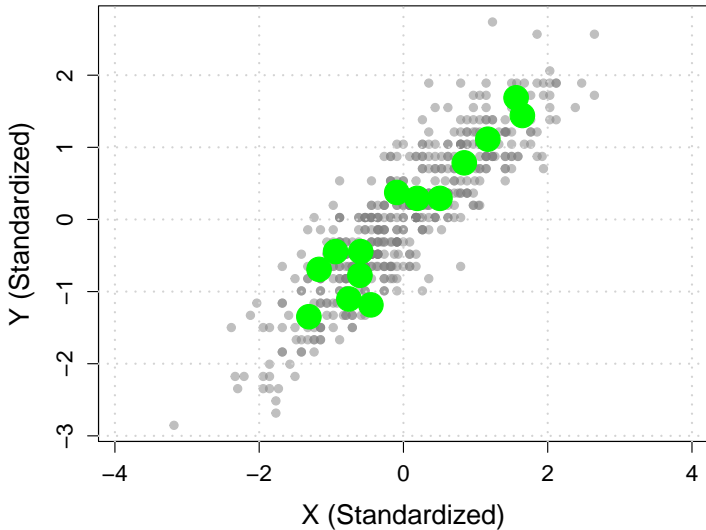


$$r = 0.8828$$

**Pups**



## Standardized scatterplots



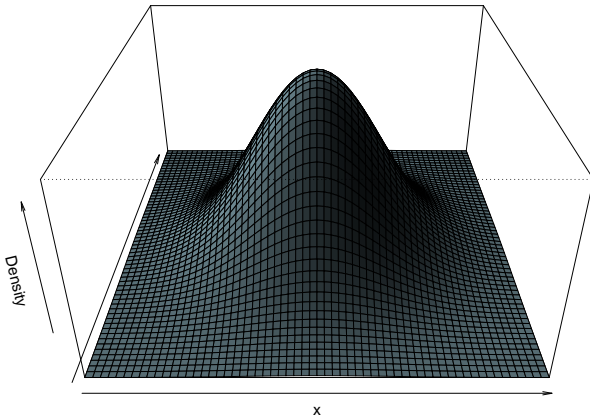


## Summary statistics and Parameters

| Summary statistic       |           | Parameter |          |
|-------------------------|-----------|-----------|----------|
| sample mean             | $\bar{x}$ | mean      | $\mu$    |
| sample s.d.             | $s_x$     | s.d.      | $\sigma$ |
| correlation coefficient | $r$       | corr.     | $\rho$   |

## Bivariate normal distribution

$$f(x, y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left[ \frac{(x-\mu_x)^2}{\sigma_x^2} + \frac{(y-\mu_y)^2}{\sigma_y^2} - \frac{2\rho(x-\mu_x)(y-\mu_y)}{\sigma_x\sigma_y} \right]}$$



## Bivariate normal distribution

$$f(x, y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left[ \frac{(x-\mu_x)^2}{\sigma_x^2} + \frac{(y-\mu_y)^2}{\sigma_y^2} - \frac{2\rho(x-\mu_x)(y-\mu_y)}{\sigma_x\sigma_y} \right]}$$

Note:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1$$

And:

$$\Pr(x < A \text{ and } y < B) = \int_{-\infty}^A \int_{-\infty}^B f(x, y) dx dy$$

## Bivariate normal distribution

$$f(x, y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left[ \frac{(x-\mu_x)^2}{\sigma_x^2} + \frac{(y-\mu_y)^2}{\sigma_y^2} - \frac{2\rho(x-\mu_x)(y-\mu_y)}{\sigma_x\sigma_y} \right]}$$

Note:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1$$

And:

$$\Pr(x < A \text{ and } y < B) = \int_{-\infty}^A \int_{-\infty}^B f(x, y) dx dy$$

## Bivariate normal distribution

$$f(x, y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left[ \frac{(x-\mu_x)^2}{\sigma_x^2} + \frac{(y-\mu_y)^2}{\sigma_y^2} - \frac{2\rho(x-\mu_x)(y-\mu_y)}{\sigma_x\sigma_y} \right]}$$

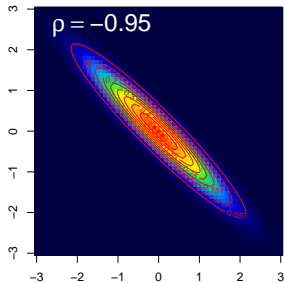
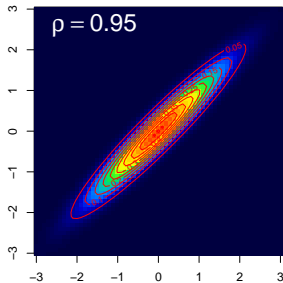
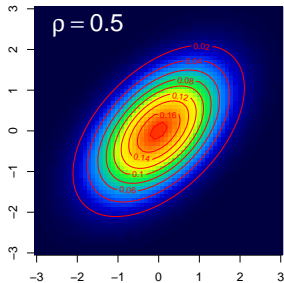
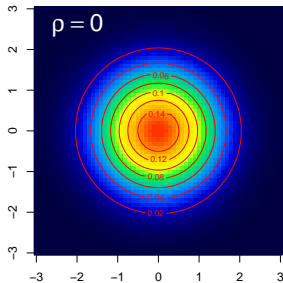
Note:

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1$$

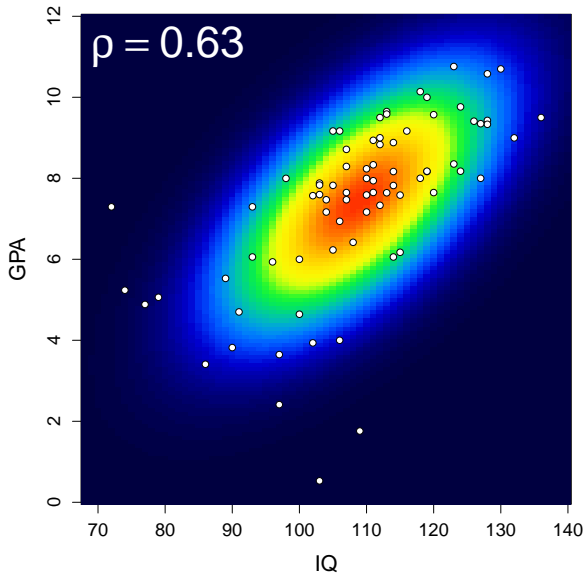
And:

$$\Pr(x < A \text{ and } y < B) = \int_{-\infty}^A \int_{-\infty}^B f(x, y) dx dy$$

# Bivariate normal distribution



## Bivariate normal distribution



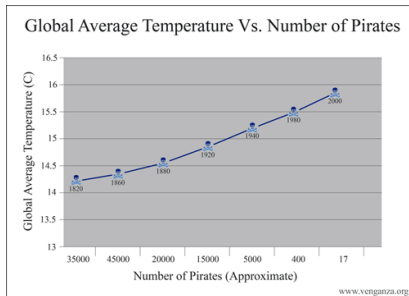
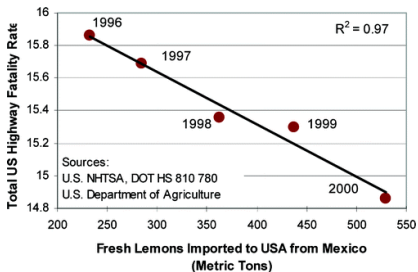
## Correlations are...

- unitless and independent of units of measurement;
- between -1 (perfect negative) and +1 (perfect positive) with  $r = 0$  meaning no relationship;
- symmetric (no separation between “cause” and “effect”).



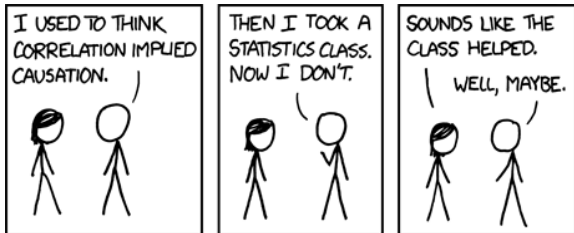
# Correlation does not imply Causation!

All calculating **correlations** does is suggest the strength and direction of the relationship between two variables.



It is easy to find numbers that are related due to **confounding** or **hidden** variable (note in these examples above the crucial hidden variable of TIME).

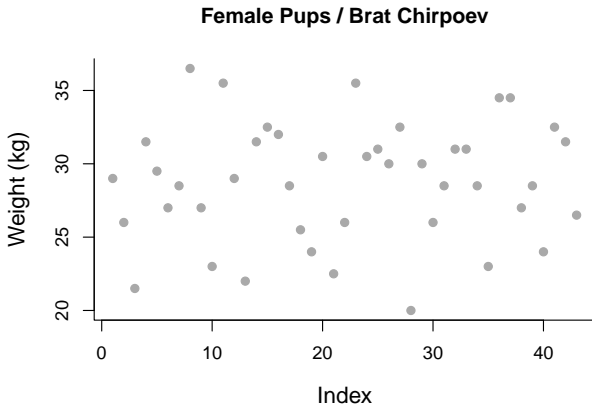
## Or does it?



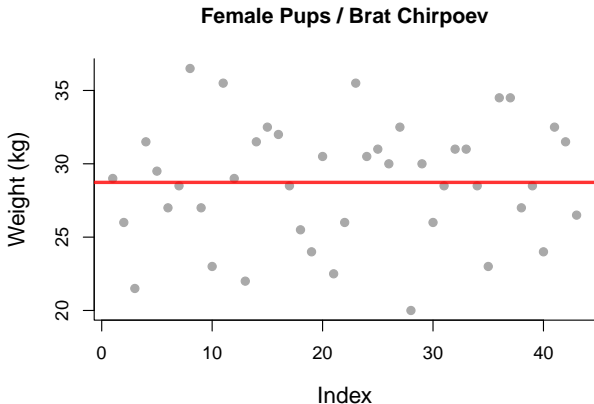
## Part II: Linear regression models

## A brief review of estimating means and s.d.'s

## A brief review of estimating means and s.d.'s

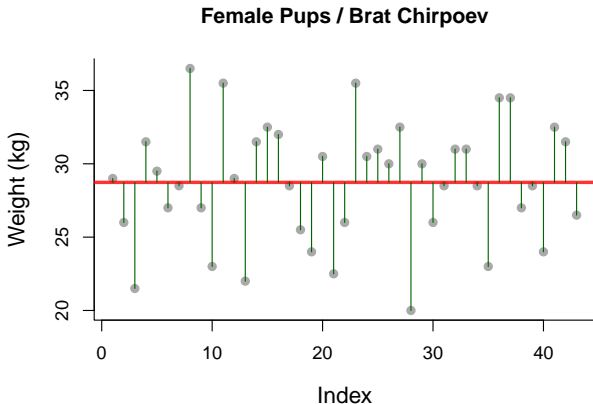


## A brief review of estimating means and s.d.'s



$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

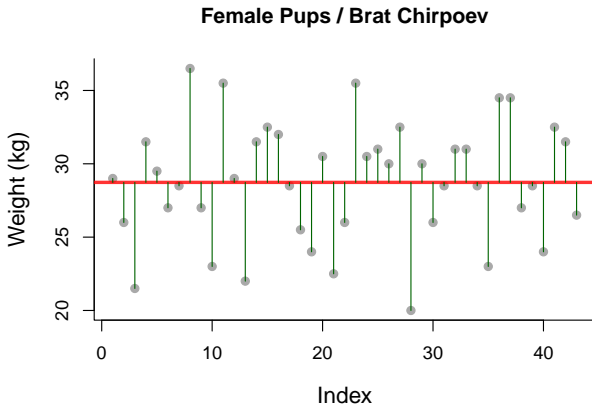
## A brief review of estimating means and s.d.'s



$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

# Formulating a model



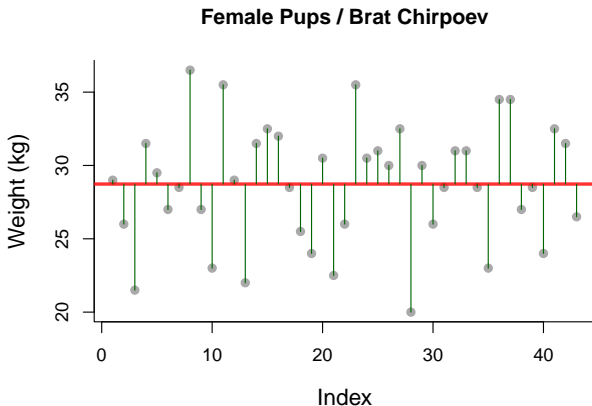
There are two ways to write this model:

$$W \sim N(\mu = \bar{X}, \sigma^2 = s_x^2)$$

$$\begin{aligned} W &= \bar{X} + \epsilon_i \\ \text{where: } \epsilon &\sim N(0, \sigma^2 = s_x^2) \end{aligned}$$



# Formulating a model



There are two ways to write this model:

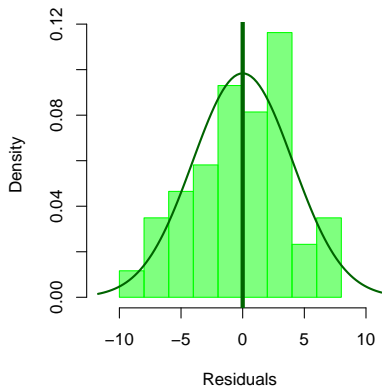
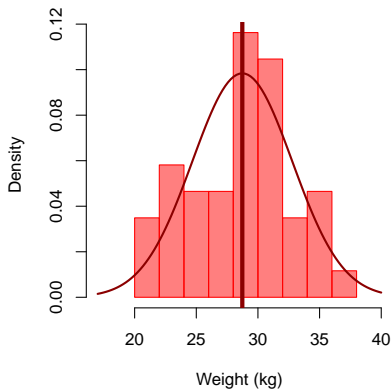
$$W \sim N(\mu = \bar{X}, \sigma^2 = s_x^2)$$

$$W = \bar{X} + \epsilon_i$$

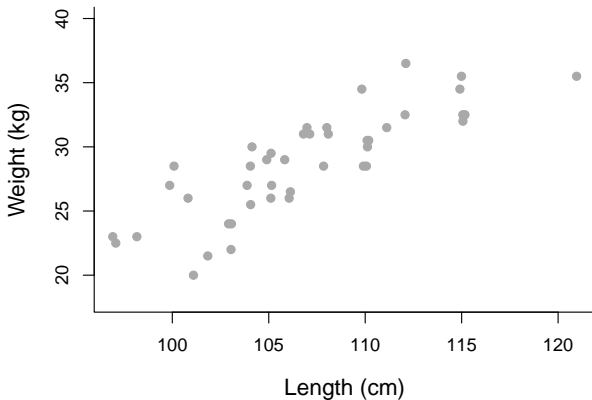
where:  $\epsilon \sim N(0, \sigma^2 = s_x^2)$

$\epsilon$ 's are called the **deviations** or the **residuals**

## Histogram of residuals

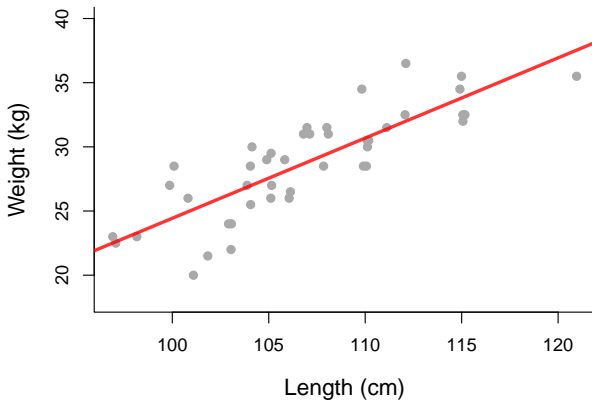


## What if two variables are related?



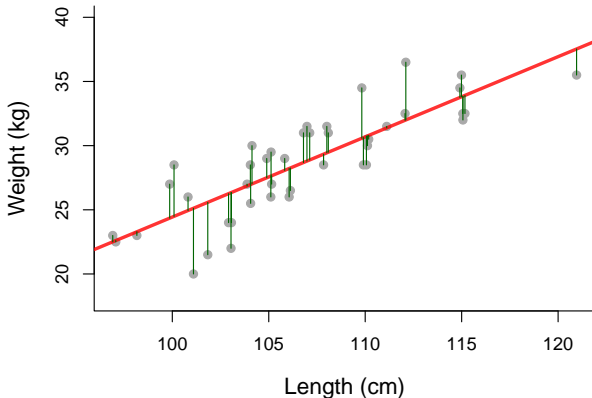
- Step 1: Draw the points

## What if two variables are related?



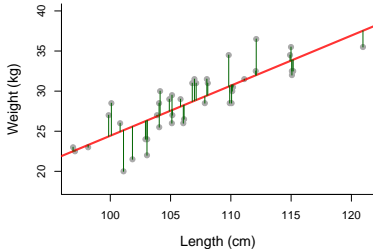
- Step 1: Draw the points
- Step 2: Write a model:  $Y_i = \alpha + \beta X_i + \epsilon$

## What if two variables are related?



- Step 1: Draw the points
- Step 2: Write a model:  $Y_i = \alpha + \beta X_i + \epsilon$
- Step 3: Calculate residuals.

# The Model



Linear model:  $Y_i = \alpha + \beta X_i + \epsilon$ .

$\alpha$  is the **intercept**

- tells us what  $Y$  would be if  $X$  were 0.
- units: same as  $Y$

$\beta$  is the **slope**

- tells how much  $Y$  will increase with each increment of  $X$
- units:  $Y$ -units/ $X$ -units

$\epsilon$  are **residuals**

- A possible (common) *model* for residuals is i.i.d.  $N(0, \sigma^2)$

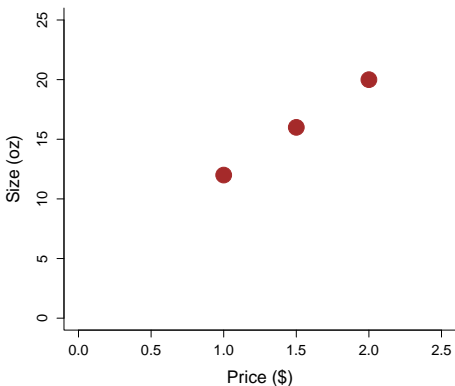
## A review of lines

|                 | Price (\$) | Size (oz) |
|-----------------|------------|-----------|
| Tall (small)    | 1.00       | 12        |
| Grande (medium) | 1.50       | 16        |
| Vente (large)   | 2.00       | 20        |



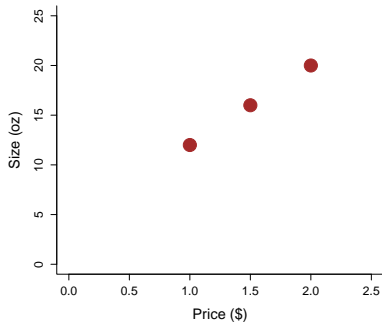
## Coffee scatterplot

|                 | Price (\$) | Size (oz) |
|-----------------|------------|-----------|
| Tall (small)    | 1.00       | 12        |
| Grande (medium) | 1.50       | 16        |
| Vente (large)   | 2.00       | 20        |





# Coffee scatterplot



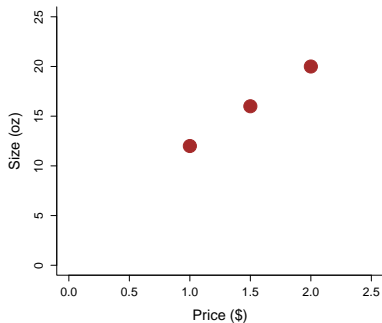
$$y = a + bx$$

$$\begin{aligned} b &= \frac{\Delta y}{\Delta x} = \frac{y_3 - y_1}{x_3 - x_1} \\ &= 8/1 = 8 \text{ oz}/\$ \end{aligned}$$

$$\begin{aligned} a &= y_1 - b x_1 \\ &= 12 - 8 \times 1 = 4 \text{ oz}. \end{aligned}$$

$$y = 4 + 8x$$

# Coffee scatterplot



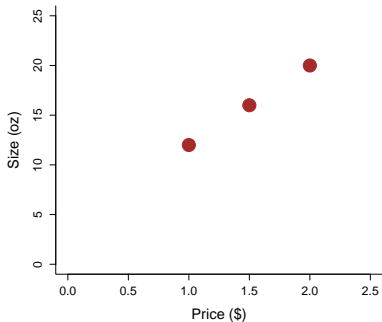
$$y = a + bx$$

$$\begin{aligned} b &= \frac{\Delta y}{\Delta x} = \frac{y_3 - y_1}{x_3 - x_1} \\ &= 8/1 = 8 \text{ oz}/\$ \end{aligned}$$

$$\begin{aligned} a &= y_1 - b x_1 \\ &= 12 - 8 \times 1 = 4 \text{ oz}. \end{aligned}$$

$$y = 4 + 8x$$

# Coffee scatterplot



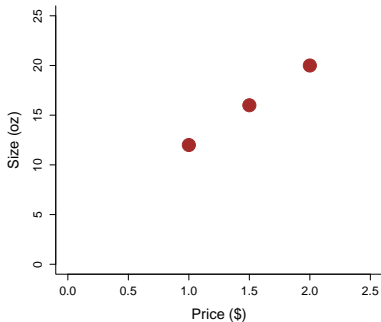
$$y = a + bx$$

$$\begin{aligned} b &= \frac{\Delta y}{\Delta x} = \frac{y_3 - y_1}{x_3 - x_1} \\ &= 8/1 = 8 \text{ oz}/\$ \end{aligned}$$

$$\begin{aligned} a &= y_1 - b x_1 \\ &= 12 - 8 \times 1 = 4 \text{ oz}. \end{aligned}$$

$$y = 4 + 8x$$

# Coffee scatterplot



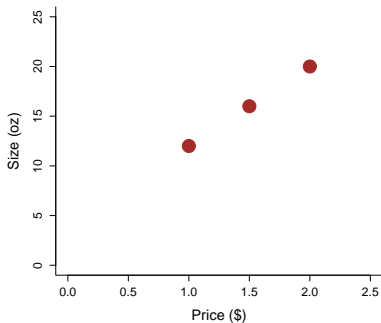
$$y = a + bx$$

$$\begin{aligned} b &= \frac{\Delta y}{\Delta x} = \frac{y_3 - y_1}{x_3 - x_1} \\ &= 8/1 = 8 \text{ oz}/\$ \end{aligned}$$

$$\begin{aligned} a &= y_1 - bx_1 \\ &= 12 - 8 \times 1 = 4 \text{ oz}. \end{aligned}$$

$$y = 4 + 8x$$

# Coffee scatterplot



$$y = a + bx$$

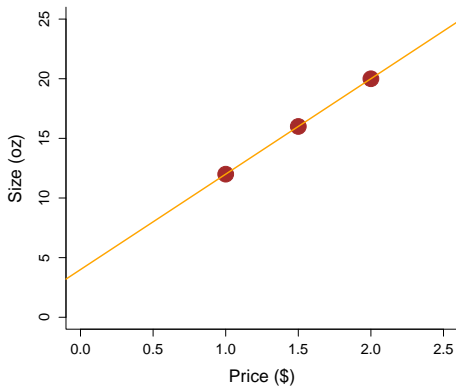
$$\begin{aligned} b &= \frac{\Delta y}{\Delta x} = \frac{y_3 - y_1}{x_3 - x_1} \\ &= 8/1 = 8 \text{ oz}/\$ \end{aligned}$$

$$\begin{aligned} a &= y_1 - bx_1 \\ &= 12 - 8 \times 1 = 4 \text{ oz}. \end{aligned}$$

$$y = 4 + 8x$$

## Coffee scatterplot

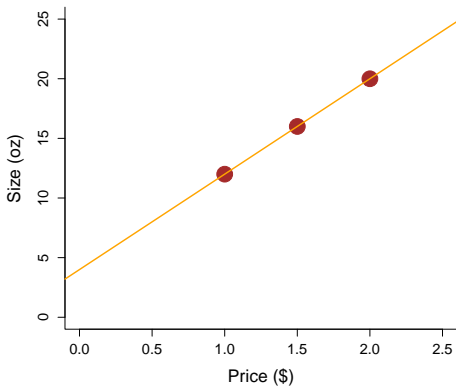
$$y = 4 + 8x$$



We have proven that: a 4 oz. coffee costs \$0!

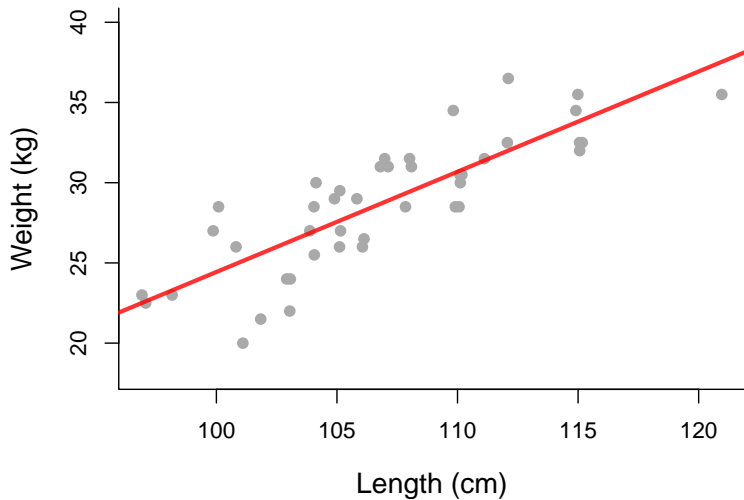
## Coffee scatterplot

$$y = 4 + 8x$$



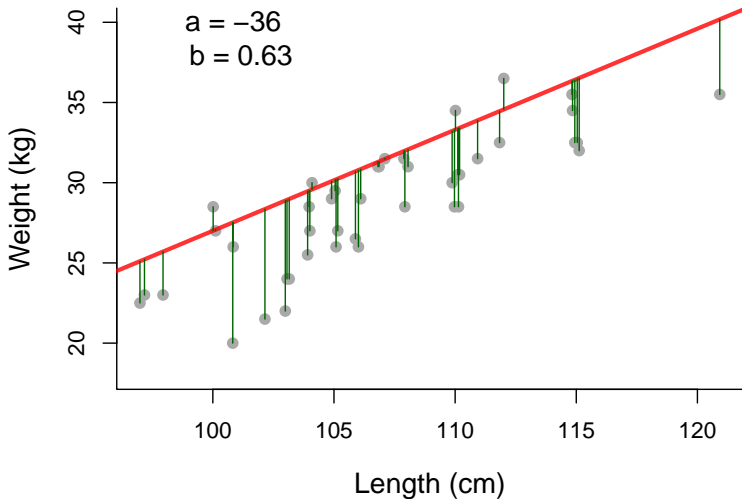
We have proven that: a 4 oz. coffee costs \$0!

But how do we pick the line if the points are scattered?

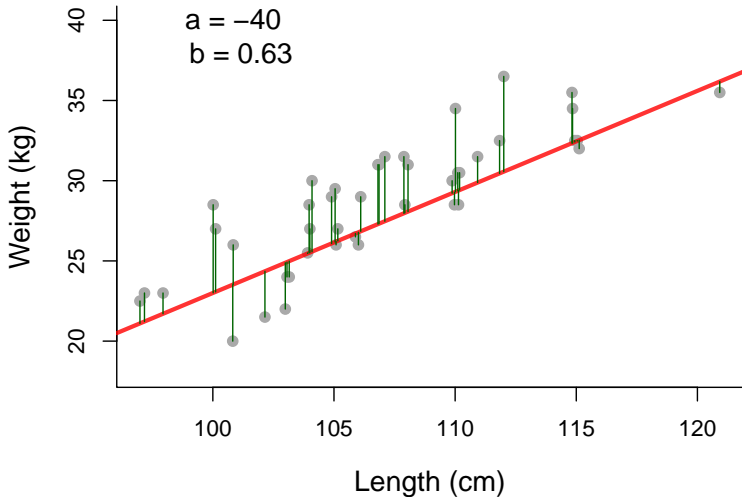




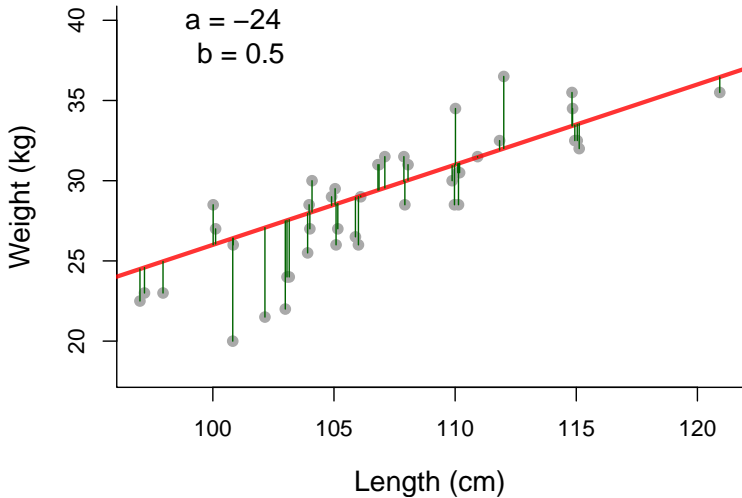
Lots and lots of lines are possible!



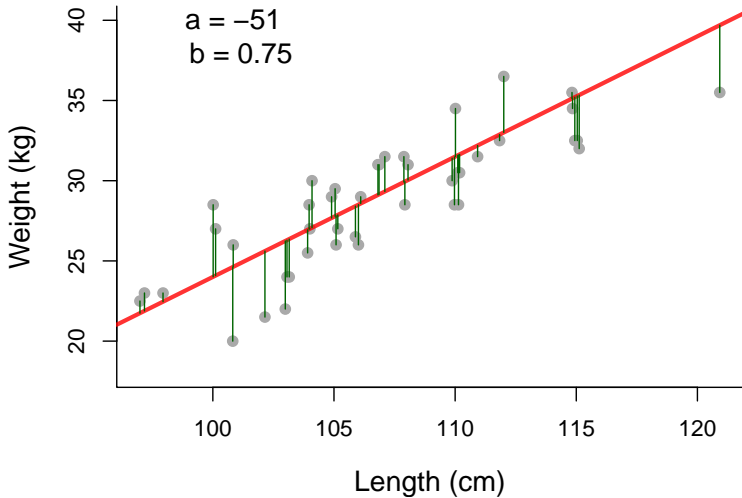
Lots and lots of lines are possible!



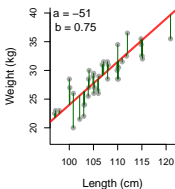
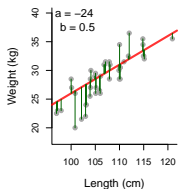
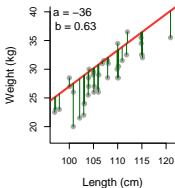
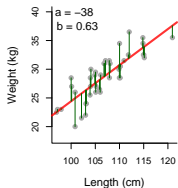
Lots and lots of lines are possible!



Lots and lots of lines are possible!

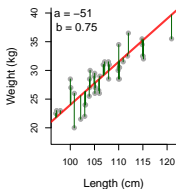
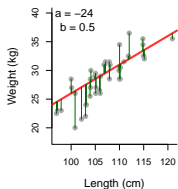
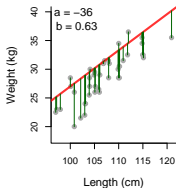
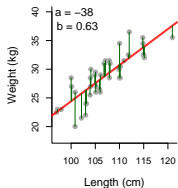


# How do we know it is a good line?



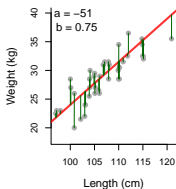
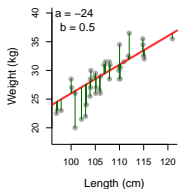
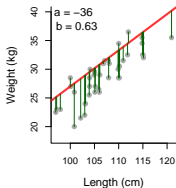
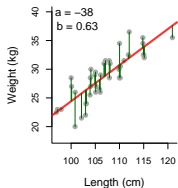
- Calculate residuals:  
 $\epsilon_i = Y_i - (a + bX_i);$
- Sum their squares ( $SS_{error}$ )  
 $SS_{error} = \sum \epsilon_i^2 = \sum_{i=1}^n (Y_i - (a + bX_i))^2;$
- Find values of  $a$  and  $b$  that minimize the  $SS_{error}$ ;

# How do we know it is a good line?



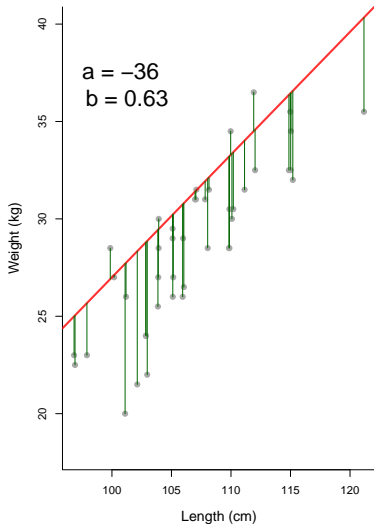
- Calculate residuals:  
 $\epsilon_i = Y_i - (a + bX_i);$
- Sum their squares ( $SS_{error}$ )  
$$SS_{error} = \sum \epsilon_i = \sum_{i=1}^n (Y_i - (a + bX_i))^2;$$
- Find values of  $a$  and  $b$  that minimize the  $SS_{error}$ ;

# How do we know it is a good line?



- Calculate residuals:  
$$\epsilon_i = Y_i - (a + bX_i);$$
- Sum their squares ( $SS_{error}$ )  
$$SS_{error} = \sum \epsilon_i = \sum_{i=1}^n (Y_i - (a + bX_i))^2;$$
- Find values of  $a$  and  $b$  that minimize the  $SS_{error}$ ;

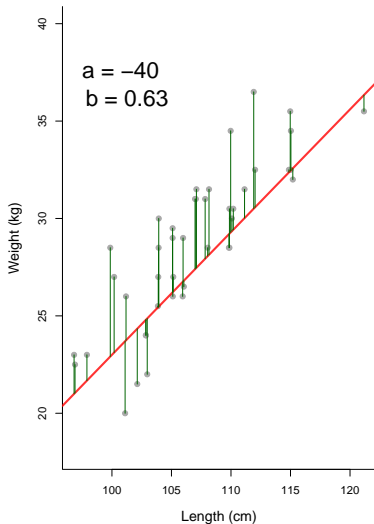
## Calculation of $SS_{error}$



| W (data)     | W (model) | Residuals       |
|--------------|-----------|-----------------|
| 29           | 30.78     | -1.78           |
| 26           | 27.63     | -1.63           |
| 21.5         | 28.26     | -6.76           |
| 31.5         | 32.04     | -0.54           |
| 29.5         | 30.15     | -0.65           |
| 27           | 27        | 0               |
| 28.5         | 33.3      | -4.8            |
| 36.5         | 34.56     | 1.94            |
| 27           | 30.15     | -3.15           |
| 23           | 25.11     | -2.11           |
| ...          |           |                 |
| $SS_{error}$ |           | <b>500.2329</b> |

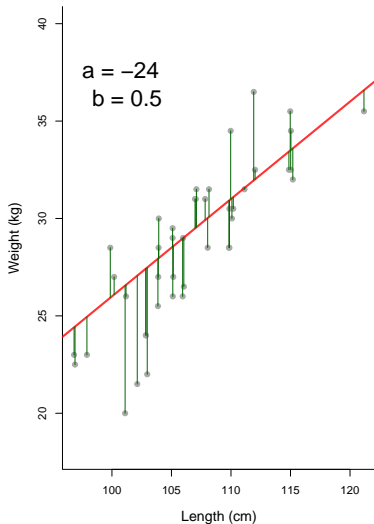


## Calculation of $SS_{error}$



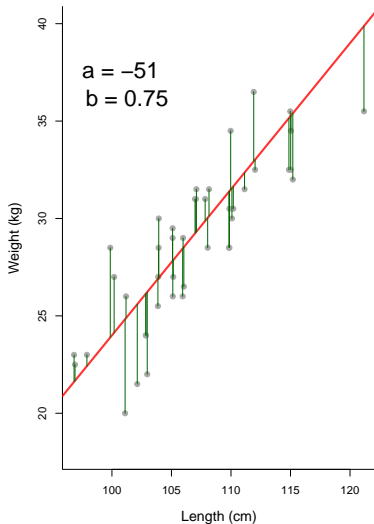
| W (data)     | W (model) | Residuals       |
|--------------|-----------|-----------------|
| 29           | 26.78     | 2.22            |
| 26           | 23.63     | 2.37            |
| 21.5         | 24.26     | -2.76           |
| 31.5         | 28.04     | 3.46            |
| 29.5         | 26.15     | 3.35            |
| 27           | 23        | 4               |
| 28.5         | 29.3      | -0.8            |
| 36.5         | 30.56     | 5.94            |
| 27           | 26.15     | 0.85            |
| 23           | 21.11     | 1.89            |
| ...          |           |                 |
| $SS_{error}$ |           | <b>297.4329</b> |

## Calculation of $SS_{error}$



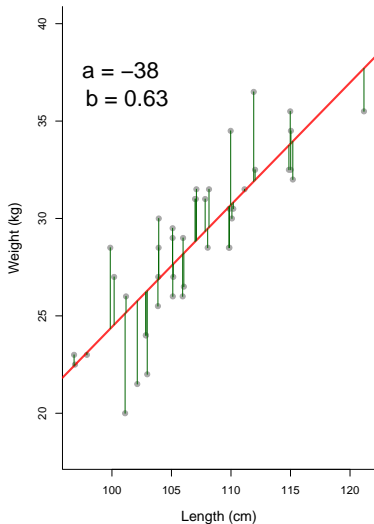
| W (data)     | W (model) | Residuals    |
|--------------|-----------|--------------|
| 29           | 29        | 0            |
| 26           | 26.5      | -0.5         |
| 21.5         | 27        | -5.5         |
| 31.5         | 30        | 1.5          |
| 29.5         | 28.5      | 1            |
| 27           | 26        | 1            |
| 28.5         | 31        | -2.5         |
| 36.5         | 32        | 4.5          |
| 27           | 28.5      | -1.5         |
| 23           | 24.5      | -1.5         |
| ...          |           |              |
| $SS_{error}$ |           | <b>252.5</b> |

## Calculation of $SS_{error}$



| W (data)     | W (model) | Residuals       |
|--------------|-----------|-----------------|
| 29           | 28.5      | 0.5             |
| 26           | 24.75     | 1.25            |
| 21.5         | 25.5      | -4              |
| 31.5         | 30        | 1.5             |
| 29.5         | 27.75     | 1.75            |
| 27           | 24        | 3               |
| 28.5         | 31.5      | -3              |
| 36.5         | 33        | 3.5             |
| 27           | 27.75     | -0.75           |
| 23           | 21.75     | 1.25            |
| ...          |           |                 |
| $SS_{error}$ |           | <b>237.5625</b> |

## Calculation of $SS_{error}$



| W (data)     | W (model) | Residuals    |
|--------------|-----------|--------------|
| 29           | 28.19     | 0.81         |
| 26           | 25.06     | 0.94         |
| 21.5         | 25.68     | -4.18        |
| 31.5         | 29.45     | 2.05         |
| 29.5         | 27.57     | 1.93         |
| 27           | 24.43     | 2.57         |
| 28.5         | 30.7      | -2.2         |
| 36.5         | 31.96     | 4.54         |
| 27           | 27.57     | -0.57        |
| 23           | 22.55     | 0.45         |
| ...          |           |              |
| $SS_{error}$ |           | <b>211.9</b> |

## How do we find the optimal $a$ and $b$ ?

- Do a lot of guessing and checking.
- Ask the computer.
- Do some fun calculus!

## How do we find the optimal $a$ and $b$ ?

- Do a lot of guessing and checking.
- Ask the computer.
- Do some fun calculus!

## How do we find the optimal $a$ and $b$ ?

- Do a lot of guessing and checking.
- Ask the computer.
- Do some fun calculus!

## Minimizing the $SS_{error}$

- Note that  $SS_{error} = f(a, b|X, Y)$ 
  - the vertical bar “|” means: “given” or **conditional**
  - so the eq. above reads - “ $SS_{error}$  is a function of parameters  $a$  and  $b$  given a known set of data  $X$  and  $Y$ ”

$$\frac{\partial f(a, b|X, Y)}{\partial a} = \frac{\partial}{\partial a} \left( \sum_{i=1}^n (Y_i - (a + bX_i))^2 \right) \equiv 0$$

$$\frac{\partial f(a, b|X, Y)}{\partial b} = \frac{\partial}{\partial b} \left( \sum_{i=1}^n (Y_i - (a + bX_i))^2 \right) \equiv 0$$

Recall that the MINIMUM occurs where the SLOPE of a function is 0, and that the DERIVATIVE tells you the SLOPE.



## Minimizing the $SS_{error}$

- Note that  $SS_{error} = f(a, b|X, Y)$ 
  - the vertical bar “|” means: “given” or **conditional**
  - so the eq. above reads - “ $SS_{error}$  is a function of parameters  $a$  and  $b$  given a known set of data  $X$  and  $Y$ ”

$$\frac{\partial f(a, b|X, Y)}{\partial a} = \frac{\partial}{\partial a} \left( \sum_{i=1}^n (Y_i - (a + bX_i))^2 \right) \equiv 0$$

$$\frac{\partial f(a, b|X, Y)}{\partial b} = \frac{\partial}{\partial b} \left( \sum_{i=1}^n (Y_i - (a + bX_i))^2 \right) \equiv 0$$

Recall that the MINIMUM occurs where the SLOPE of a function is 0, and that the DERIVATIVE tells you the SLOPE.

## Solving for the intercept

$$\frac{\partial f(a, b|X, Y)}{\partial a} = 2 \sum_{i=1}^n (Y_i - (a + bX_i)) = 0$$

Recall:  $\sum_{i=1}^n Y_i = n\bar{Y}$  and  $\sum_{i=1}^n X_i = n\bar{X}$  and  $\sum_{i=1}^n a = na$ .

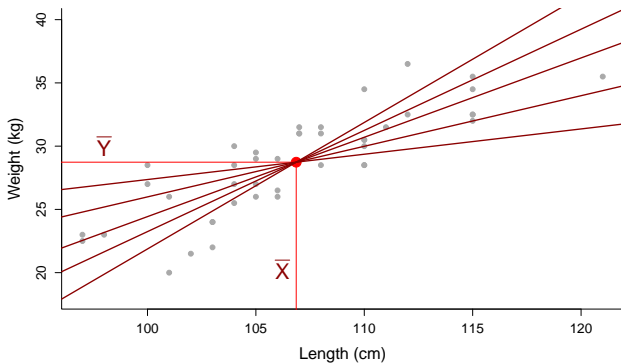
$$\begin{aligned} \sum_{i=1}^n (Y_i - (a + bX_i)) &= 0 \\ n\bar{Y} - na - nb\bar{X} &= 0 \end{aligned}$$

Leads to:

$$a = \bar{Y} - b\bar{X}$$

## Solving for the intercept

$$a = \bar{Y} - b\bar{X}$$



Implies that the regression line goes through  $\bar{X}$  and  $\bar{Y}$ .

## Solving for the slope

$$\frac{\partial f(a, b | X, Y)}{\partial b} = 2 \sum_{i=1}^n (Y_i - (a + bX_i))X_i = 0$$

Plug in  $a = \bar{Y} - b\bar{X}$

$$\sum_{i=1}^n (Y_i - (\bar{Y} - b\bar{X} + bX_i))X_i = 0$$

$$\sum (Y_i X_i - \bar{Y} X_i) - b \sum (X_i^2 - \bar{X} X_i) = 0$$

Leads to:

$$b = \frac{\sum (Y_i X_i - \bar{Y} X_i)}{\sum (X_i^2 - \bar{X} X_i)}$$

$$b = \frac{\sum(Y_i X_i - \bar{Y} X_i)}{\sum(X_i^2 - \bar{X} X_i)}$$

Note the following identities:

$$\begin{aligned}\sum \bar{X} Y_i &= \sum X_i \bar{Y} = \sum \bar{X} \bar{Y} \\ \sum X_i \bar{X} &= \sum \bar{X}^2\end{aligned}$$

Rewrite (with some algebra):

$$b = \frac{\sum(Y_i X_i - \bar{Y} X_i - \bar{X} Y_i + \bar{X} \bar{Y})}{\sum(X_i^2 - 2X_i \bar{X} + \bar{X}^2)}$$

and format:

$$b = \frac{\sum(Y_i - \bar{Y})(X_i - \bar{X})}{\sum(X_i - \bar{X})^2}$$

# Intercept and slope

**Slope:**

$$b = \frac{\sum(Y_i - \bar{Y})(X_i - \bar{X})}{\sum(X_i - \bar{X})^2}$$

**Intercept:**

$$a = \bar{Y} - b\bar{X}$$

(plug in the right value for  $b$ ).

## Intercept and slope

$$b = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$a = \bar{Y} - b\bar{X}$$

Recall:

$$s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

$$r_{xy} = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{Y_i - \bar{Y}}{s_y} \right) \left( \frac{X_i - \bar{X}}{s_x} \right)$$

So (after some algebra) we can rewrite  $a$  and  $b$  as:

$$b = r_{xy} \left( \frac{s_y}{s_x} \right)$$

$$a = \bar{Y} - r_{xy} \left( \frac{s_y}{s_x} \right) \bar{X}$$

# Linear regression: Pup Example

## Summary statistics:

$$\bar{x} = 106.9; s_x = 5.38$$

$$\bar{y} = 28.7; s_y = 4.06$$

$$r_{xy} = 0.83$$

## Regression coefficients:

$$b = r(s_y/s_x)$$

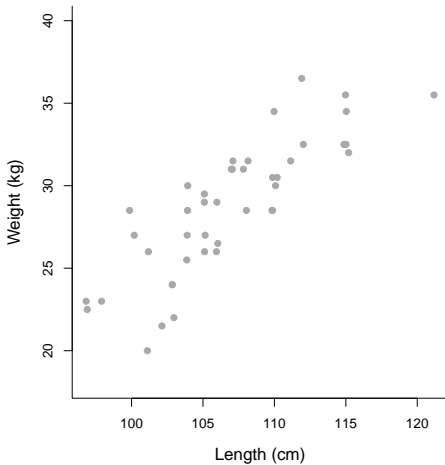
$$0.83 \times (4.06/5.38)$$

$$0.63 \text{ kg/cm}$$

$$a = \bar{y} - b\bar{x}$$

$$28.7 - 0.63 \times 106.9$$

$$-38.2 \text{ kg}$$





# Linear regression: Pup Example

## Summary statistics:

$$\bar{x} = 106.9; s_x = 5.38$$

$$\bar{y} = 28.7; s_y = 4.06$$

$$r_{xy} = 0.83$$

## Regression coefficients:

$$b = r(s_y/s_x)$$

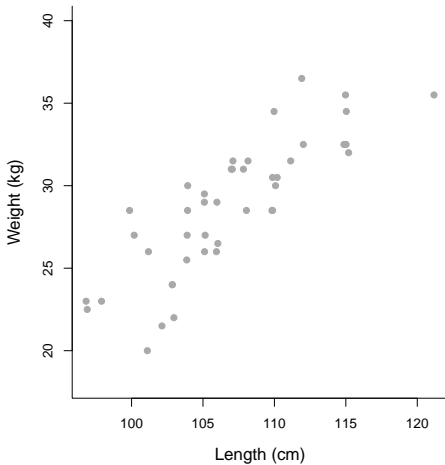
$$0.83 \times (4.06/5.38)$$

$$0.63 \text{ kg/cm}$$

$$a = \bar{y} - b\bar{x}$$

$$28.7 - 0.63 \times 106.9$$

$$-38.2 \text{ kg}$$



# Linear regression: Pup Example

## Summary statistics:

$$\bar{x} = 106.9; s_x = 5.38$$

$$\bar{y} = 28.7; s_y = 4.06$$

$$r_{xy} = 0.83$$

## Regression coefficients:

$$b = r(s_y/s_x)$$

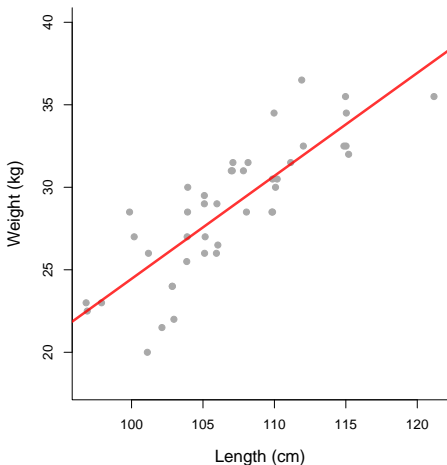
$$0.83 \times (4.06/5.38)$$

$$\mathbf{0.63 \text{ kg/cm}}$$

$$a = \bar{Y} - b\bar{X}$$

$$28.7 - 0.63 \times 106.9$$

$$\mathbf{-38.2 \text{ kg}}$$



## Important features of $\hat{Y} = a + bx$

The least squares estimates define a line with the following properties:

- The line passes through  $(\bar{X}, \bar{Y})$
- The residuals from the least squares fitted line sum to zero:

$$\sum_{i=1}^n (Y_i - \hat{Y}) = 0$$

Recalling the model:  $Y_i = a + bX_i + \epsilon$

- $\epsilon$  is distributed normally with mean 0 and estimated variance

$$\hat{\sigma}_{error}^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \hat{Y})^2$$

## Important features of $\hat{Y} = a + bx$

The least squares estimates define a line with the following properties:

- The line passes through  $(\bar{X}, \bar{Y})$
- The residuals from the least squares fitted line sum to zero:

$$\sum_{i=1}^n (Y_i - \hat{Y}) = 0$$

Recalling the model:  $Y_i = a + bX_i + \epsilon$

- $\epsilon$  is distributed normally with mean 0 and estimated variance

$$\hat{\sigma}_{error}^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \hat{Y})^2$$

## Important features of $\hat{Y} = a + bx$

The least squares estimates define a line with the following properties:

- The line passes through  $(\bar{X}, \bar{Y})$
- The residuals from the least squares fitted line sum to zero:

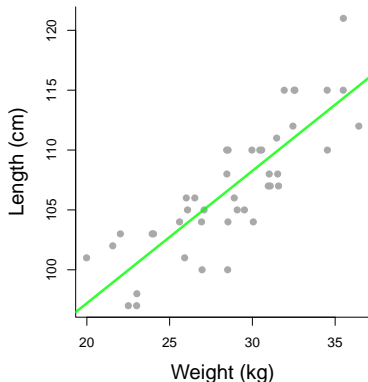
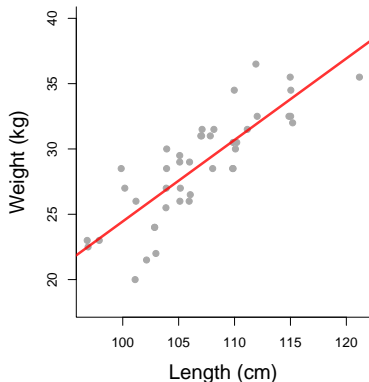
$$\sum_{i=1}^n (Y_i - \hat{Y}) = 0$$

Recalling the model:  $Y_i = a + bX_i + \epsilon$

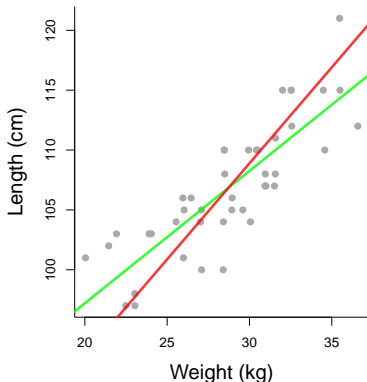
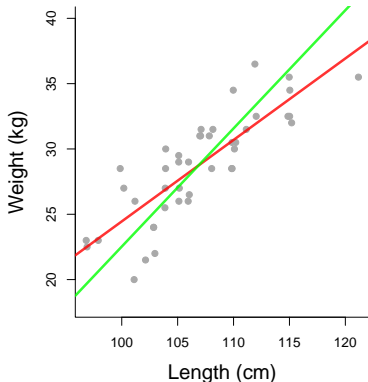
- $\epsilon$  is distributed normally with mean 0 and estimated variance

$$\hat{\sigma}_{error}^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \hat{Y})^2$$

Assymetry warning:  $b(Y|X) \neq b(X|Y)$



Assymetry warning:  $b(Y|X) \neq b(X|Y)$



But the relationship is simple:

$$b(Y|X) = r_{X,Y} \frac{s_y}{s_x} \text{ and } b(X|Y) = r_{X,Y} \frac{s_x}{s_y}$$

$$\text{so: } b(Y|X) = b(X|Y) \frac{s_y^2}{s_x^2}$$



# Some sums of squares

Sum of squares - TOTAL:

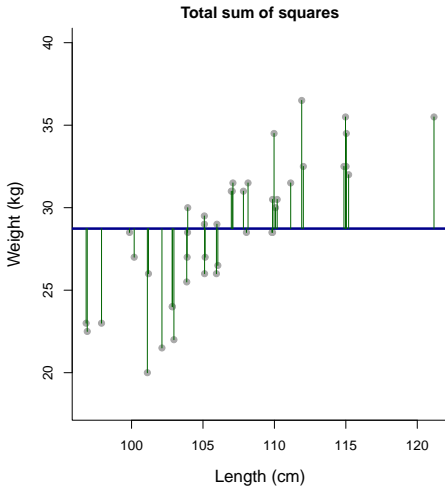
$$SS_{total} = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

MODEL sum of squares:

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

ERROR sum of squares:

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



## Some sums of squares

Sum of squares - TOTAL:

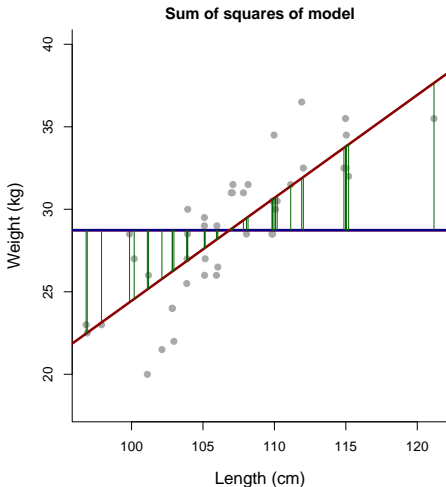
$$SS_{total} = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

MODEL sum of squares:

$$SS_{model} = \sum_{i=1}^n (\hat{Y} - \bar{Y})^2$$

ERROR sum of squares:

$$SS_{error} = \sum_{i=1}^n (Y_i - \hat{Y})^2$$



# Some sums of squares

Sum of squares - TOTAL:

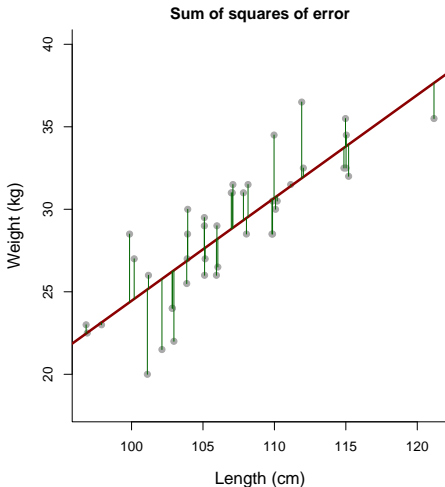
$$SS_{total} = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

MODEL sum of squares:

$$SS_{model} = \sum_{i=1}^n (\hat{Y} - \bar{Y})^2$$

ERROR sum of squares:

$$SS_{error} = \sum_{i=1}^n (Y_i - \hat{Y})^2$$



## Decomposing the total variance

- We **decompose** the total variation into “explained” and “unexplained” components. So:

Total sum of squares = Regression sum of squares + Error sum of squares

$$SS_{total} = SS_{model} + SS_{error}$$

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

## $r^2$ - coefficient of determination

$r^2$  is the proportion of total variance explained.

$$\begin{aligned} r^2 &= \frac{SS_{total} - SS_{error}}{SS_{total}} = \frac{SS_{model}}{SS_{total}} \\ &= \frac{\sum_{i=1}^n (\hat{Y} - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\ &= \frac{\sum_{i=1}^n (a + bX_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \end{aligned}$$

after plugging in  $a$  and  $b$  and lots of (not so fun) algebra

$$r^2 = \frac{\left( \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \right)^2}{s_x^2 s_y^2} = (r)^2$$

$r^2$  is a summary statistic that measures the proportion of variability explained by the model. In linear regression (but not in general)  $r^2$  is the coefficient of correlation squared.

## $r^2$ - coefficient of determination

$r^2$  is the proportion of total variance explained.

$$\begin{aligned} r^2 &= \frac{SS_{total} - SS_{error}}{SS_{total}} = \frac{SS_{model}}{SS_{total}} \\ &= \frac{\sum_{i=1}^n (\hat{Y} - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\ &= \frac{\sum_{i=1}^n (a + bX_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \end{aligned}$$

after plugging in  $a$  and  $b$  and lots of (not so fun) algebra

$$r^2 = \frac{\left( \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \right)^2}{s_x^2 s_y^2} = (r)^2$$

$r^2$  is a summary statistic that measures the proportion of variability explained by the model. In linear regression (but not in general)  $r^2$  is the coefficient of correlation squared.

## $r^2$ - coefficient of determination

$r^2$  is the proportion of total variance explained.

$$\begin{aligned} r^2 &= \frac{SS_{total} - SS_{error}}{SS_{total}} = \frac{SS_{model}}{SS_{total}} \\ &= \frac{\sum_{i=1}^n (\hat{Y} - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\ &= \frac{\sum_{i=1}^n (a + bX_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \end{aligned}$$

after plugging in  $a$  and  $b$  and lots of (not so fun) algebra

$$r^2 = \frac{\left( \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \right)^2}{s_x^2 s_y^2} = (r)^2$$

$r^2$  is a summary statistic that measures the proportion of variability explained by the model. In linear regression (but not in general)  $r^2$  is the coefficient of correlation squared.

# Linear regression: Pup Example

## Summary statistics:

$$\bar{x} = 106.9; s_x = 5.38$$

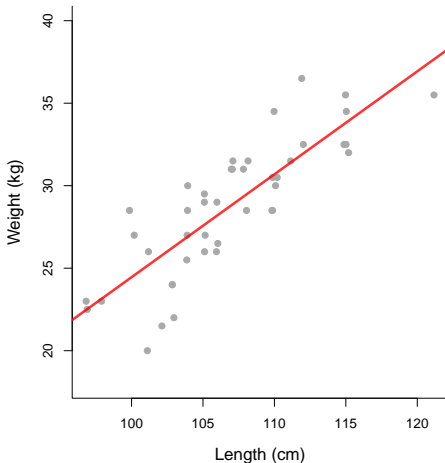
$$\bar{y} = 28.7; s_y = 4.06$$

$$r_{xy} = 0.83$$

## Coefficient of determination:

$$r^2 = 0.83^2 = 0.689$$

So we conclude: "About 70% of the observed variation in weight is explained by a linear regression against length."



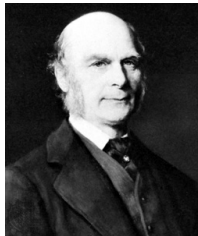


# Review of Measures of Association

| Name                         | Definition  | Comments  |
|------------------------------|---|---|
| Correlation coefficient      | $r_{xy} = \frac{1}{n-1} \sum \left( \frac{X_i - \bar{X}}{s_x} \right) \left( \frac{Y_i - \bar{Y}}{s_y} \right)$ | Unitless, Range: $(-1, 1)$<br>$r_{xy} = r_{yx}$   |
| Coefficient of determination | $r_{xy}^2 = 1 - \frac{SS_{error}}{SS_{total}} = \frac{SS_{model}}{SS_{total}}$                                  | Unitless, Range: $(0, 1)$<br>$r_{xy}^2 = r_{yx}^2$  |
| Regression coefficient       | $b(Y X) = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$                                   | Units: $u_y / u_x$<br>Range: $(-\infty, \infty)$<br>$b(Y X) \neq b(X Y)$<br>$b(Y X) = b(X Y) \frac{s_y^2}{s_x^2}$ |

# Historical roots of Linear Regression

Linear regression owes much to **Sir Francis Galton** (1822–1911), a half-cousin of Charles Darwin and one of a generation of basically brilliant English Victorian polymaths. He made important contributions to anthropology, geography, meteorology, genetics, psychometrics and statistics.



Galton was really, really into counting and quantifying things. He noted that 'exceptional' parents produce more 'mediocre' children (and, interestingly, vice versa!). Hence the idea of 'regression' (as in regression to mediocrity). This slightly misleading name has stuck to a very useful statistical tool to this day.

His contributions were truly many and diverse (note: the questionnaire! the dog whistle! forensic fingerprinting! the Galton-Watson stochastic process!) Fortunately, some of his greatest passions, *eugenics* and *phrenology*, never got too far off the ground.