



More statistics: Correlation and Regression Coefficients

Elie Gurarie

Biol 799 - Lecture 2
January 2, 2017

January 2, 2017



Correlation (r)

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

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?

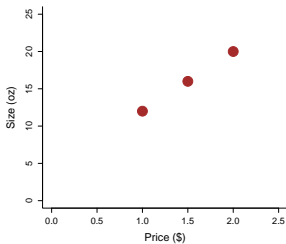
Correlation: Coffee example

	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

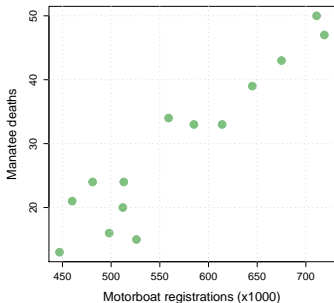
$$r = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

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

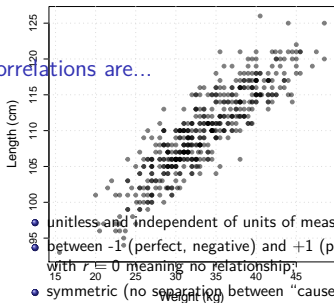
$r = 1.0$ means PERFECT correlation and a POSITIVE relationship.

More Correlations

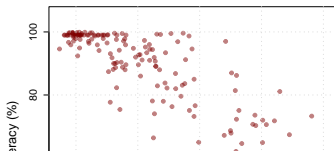
Motorboats vs. Manatees



Pups



Literacy vs. Birthrate



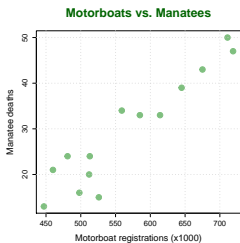
$$r = 0.9415$$

$$r = 0.8828$$

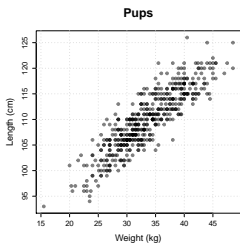
$$r = -0.8138$$

Why is...

$r = 0.9415$

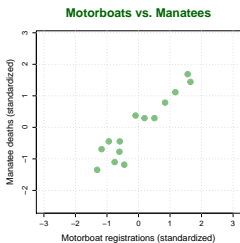


$r = 0.8828$

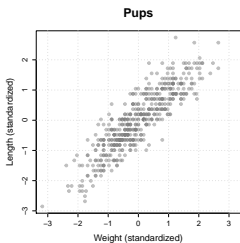


Standardized scatterplots

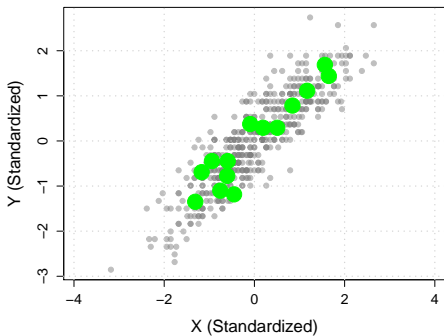
$r = 0.9415$



$r = 0.8828$



Standardized scatterplots

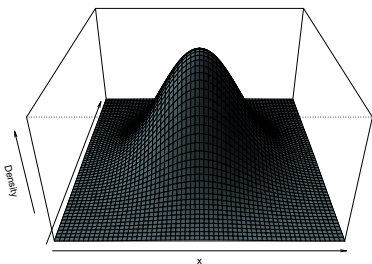


Summary statistics and Parameters

Summary statistic		Parameter	
sample mean	\bar{x}	mean	μ
sample s.d.	s_x	s.d.	σ
correlation coefficient	r	corr.	ρ

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

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

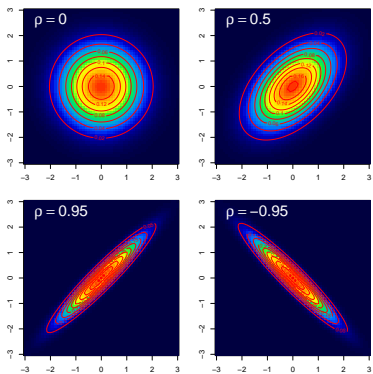
Note:

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

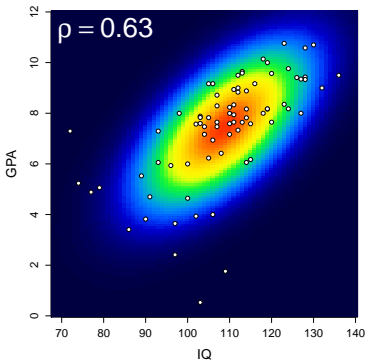
And:

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

Bivariate normal distribution



Bivariate normal distribution

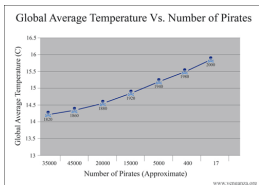
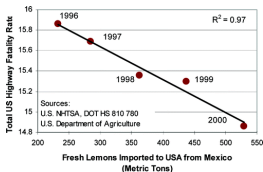


Correlations are...

- 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")**.

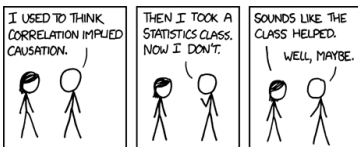
Correlation does not imply Causation!

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**).

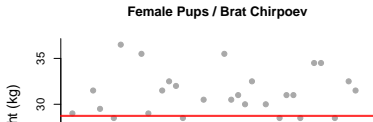
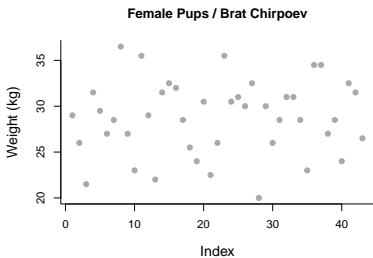
Or does it?



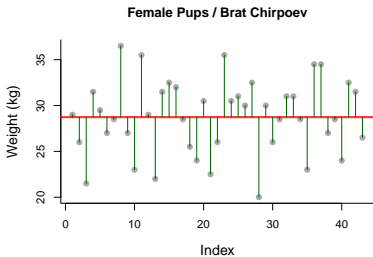
Part II: Linear regression models

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

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



Formulating a model



There are two ways to write this model:

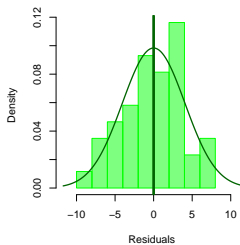
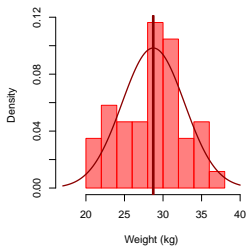
$$W \sim N(\mu = \bar{X}, \sigma^2 = s_x^2)$$

ϵ 's are called the **deviations** or the **residuals**

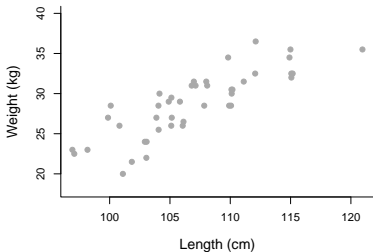
$$W = \bar{X} + \epsilon_i$$

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

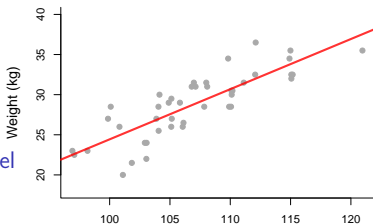
Histogram of residuals



What if two variables are related?



The Model



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

α is the **intercept**

- tells us what Y would be if X were 0.

- units: same as Y

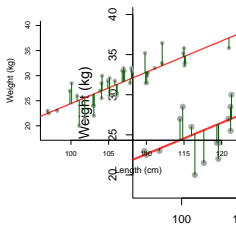
β is the **slope**

- tells how much Y will increase with each increment of X

- units: Y -units/ X -units

ϵ_i are **residuals**

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



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

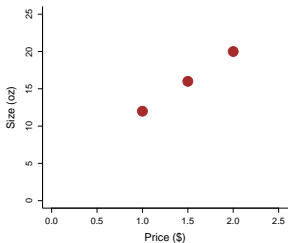
A review of lines

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

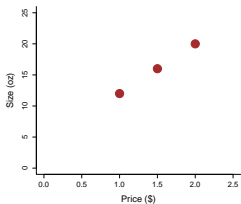


Coffee scatterplot

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



Coffee scatterplot



$$y = a + bx$$

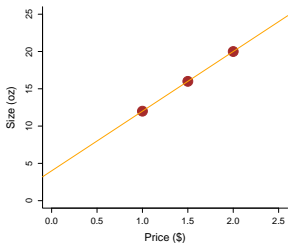
$$b = \frac{\Delta y}{\Delta x} = \frac{y_3 - y_1}{x_3 - x_1}$$
$$= \frac{8}{1} = 8 \text{ oz}/\$$$

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

$$y = 4 + 8x$$

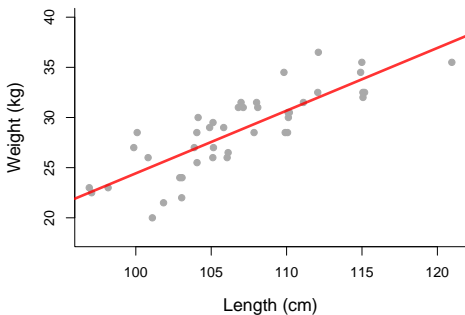
Coffee scatterplot

$$y = 4 + 8x$$

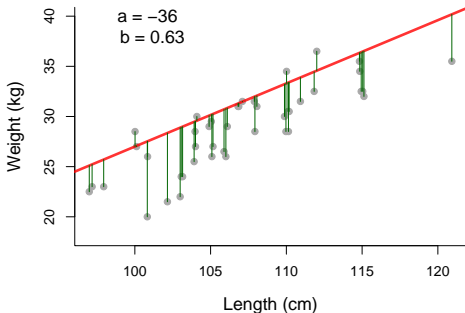


We have proven that: a 4 oz. coffee costs \$0!

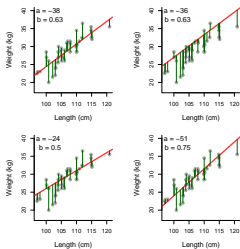
But how do we pick the line if the points are scattered?



Lots and lots of lines are possible!

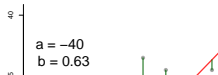
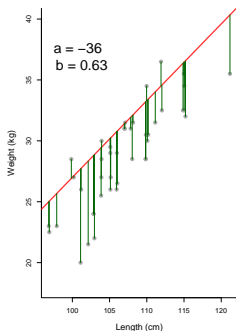


How do we know it is a good line?



- 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} ;

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

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**.

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 $\sum_{i=1}^n a = na$.

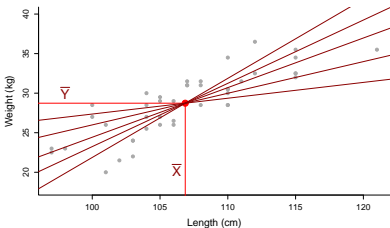
$$\begin{aligned}\sum_{i=1}^n (Y_i - (a + bX_i)) &= 0 \\ n\bar{Y} - na - nb\bar{X} &= 0\end{aligned}$$

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

$$\begin{aligned} \sum_{i=1}^n (Y_i - (\bar{Y} - b\bar{X} + bX_i))X_i &= 0 \\ \sum (Y_i X_i - \bar{Y} X_i) - b \sum (X_i^2 - \bar{X} X_i) &= 0 \end{aligned}$$

Leads to:

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

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

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:

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

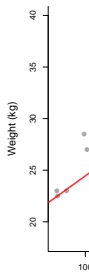
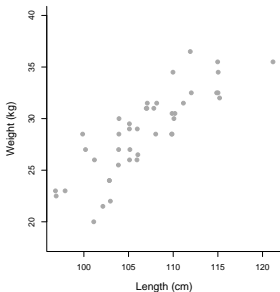
$$0.83 \times (4.06/5.38)$$

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

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

$$28.7 - 0.63 \times 106.9$$

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



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:

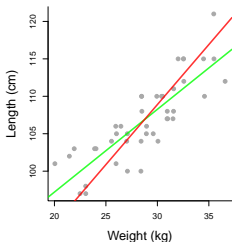
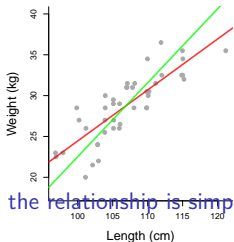
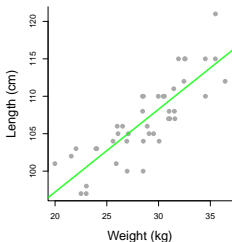
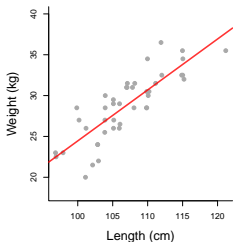
$$\sum_{i=1}^n (Y_i - \hat{Y}_i) = 0$$

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

- ϵ is distributed normally with mean 0 and estimated variance

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

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

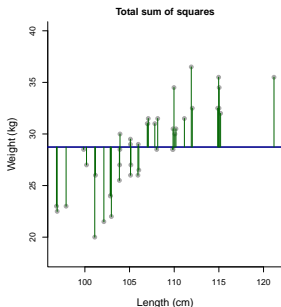


But the relationship is simple:

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

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

Some sums of squares



Sum of squares - TOTAL:

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

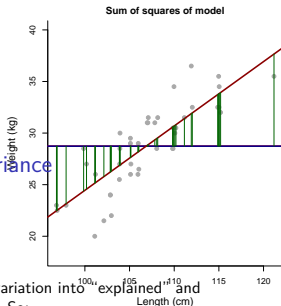
MODEL sum of squares:

Decomposing the total variance

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

ERROR sum of squares:

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

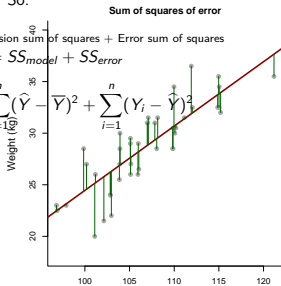


- We **decompose** the total variation into **explained** and **“unexplained”** components. So:

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

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

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



r^2 - coefficient of determination

r^2 is the proportion of total variance explained.

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

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

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

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:

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

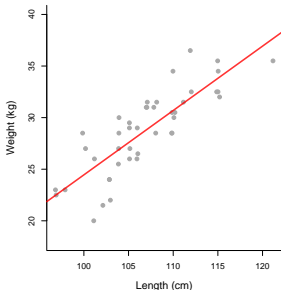
$$\bar{y} = 28.7; s_y = 4.06$$

$$r_{xy} = 0.83$$

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."

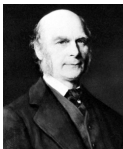


Review of Measures of Association

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

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.

Inference

We now know how to:

- Estimate means
- Estimate standard deviations
- Estimate regression coefficients

Eventually, we would like to be able to answer the following questions:

- Are means/variances of two or more samples different?
- Are regression coefficients eac
- What is the “best” model for fitting data?
- How do me make a prediction based on a fitted model?

These are all the domain of **INFERENCE!** (but first ... some Probability Theory)