

Bivariate data - Recall

- Exploratory data Analysis

Stacked
Histogram

Box plot

Data: Categorical vs Numeric

In general Bivariate data could be

Categorical vs Categorical

Numeric vs Numeric

- Setting we discussed last time

x - independent variable
Predictor or Explanatory
Variable

y - dependent variable
Response variable

- Story of least squares

Model:

y - arc length

x - function of latitude

$$y = mx + c$$

: Find m, c

from measurements $(x_1, y_1), (x_2, y_2) \dots, (x_n, y_n)$

Paired Data

$(x_i, y_i) \equiv$ ^{Dependent data} measurements from individual i

- Bivariate Data that are coupled or matched together.
They are not independent.

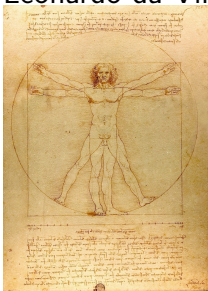
Example:

- Height and Weight measurements of individuals.
- Response reading before and after treatment of individuals.

Paired Data

Example:

- Leonardo da Vinci's *Vitruvian Man*.



- Drawing ~ 1490
- measurement of
male models

- The outstretched arms and legs within circles and square.
- Ideal human proportions described by ancient Roman architect Vitruvius: height is same as length of arm span.

Paired Data

Key Tools to understand Data

- Plot to gauge relationship.
- Correlation between the variables.
- Trends

Can they be used to
predict the %age of
body fat?

Examine ~
data set in R

- package
Using R

- fat :- Dataset

Physical
measurements of
252 males.

Paired Data

Consider `fat` dataset in `UsingR` package. The dataset contains body dimensions of 251 males.

```
> require(UsingR)
> names(fat)
```

[1]	"case"	"body.fat"	"body.fat.siri"
[4]	"density"	"age"	"weight"
[7]	"height"	"BMI"	"ffweight"
[10]	"neck"	"chest"	"abdomen"
[13]	"hip"	"thigh"	"knee"
[16]	"ankle"	"bicep"	"forearm"
[19]	"wrist"		

Paired Data

- Suppose we are interested in relation between neck and wrist.

We can first compare averages in two ways:

```
> z = mean(fat$neck)/mean(fat$wrist)
```

```
> z
```

```
[1] 2.084068
```

```
> y = mean(fat$neck/fat$wrist)
```

```
> y
```

```
[1] 2.084477
```

$$x \equiv \text{wrist}$$

$$y \equiv \text{neck}$$

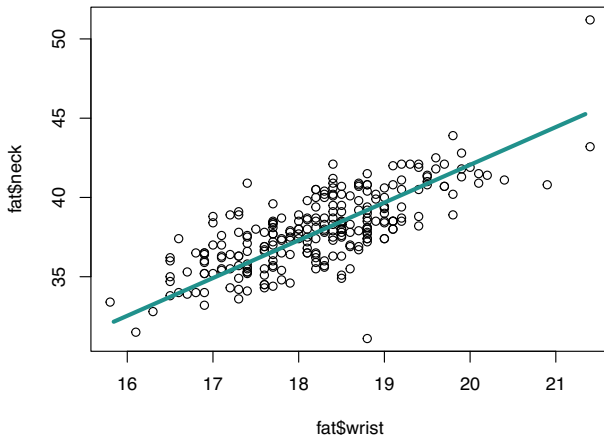
\vdots
 \downarrow

$$y = 2.08x$$

"Model"

Paired Data: dataset in UsingR

```
> plot(fat$wrist, fat$neck)
```

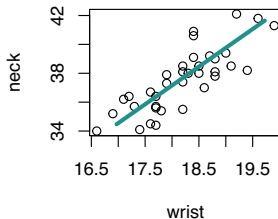
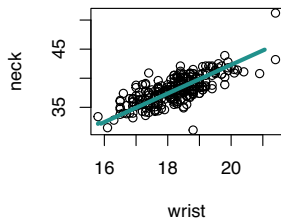


Base R

- plot
neck vs wrist
- relationship
seems linear

Paired Data: dataset in UsingR

```
> par(mfrow=c(1,2))  
> plot(neck~wrist, data=fat)  
> plot(neck~wrist, data=fat, subset=20<=age & age <30)
```



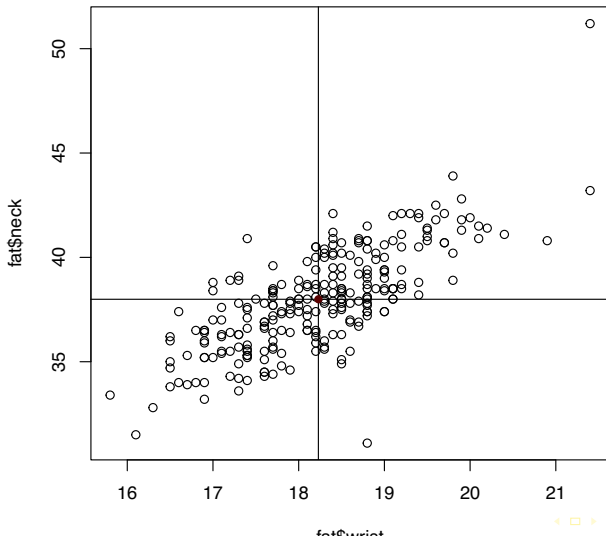
The variables seem related and also by a linear relationship

Paired Data: Correlation

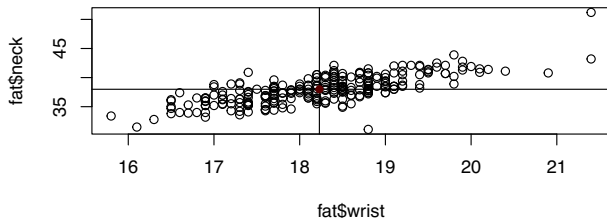
- Assume Linear Relationship between the data
- Correlation is a measure of how close the relationship is.

Before defining the term let us try to understand the plot better.

Data in four regions by means



Data in four regions by means

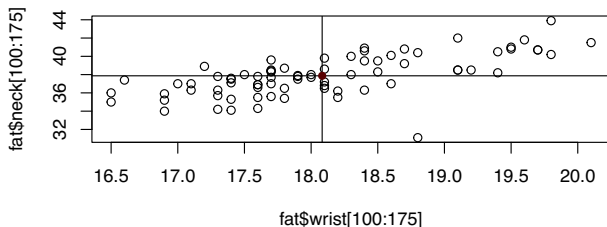


Distorted
version of
Earlier
graph

- Understand data by those above average values and those below.
- If related then most of data should be in first and third box.

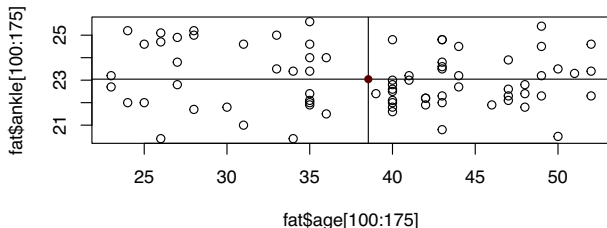
Paired Data: dataset in UsingR

```
> plot(fat$wrist[100:175], fat$neck[100:175])  
> abline(v=mean(fat$wrist[100:175]))  
> abline(h=mean(fat$neck[100:175]))  
> points(mean(fat$wrist[100:175]), mean(fat$neck[100:175]),  
+ pch=16, col=rgb(.35,0,0))
```



Paired Data: dataset in UsingR

```
> plot(fat$age[100:175], fat$ankle[100:175])  
> abline(v=mean(fat$age[100:175]))  
> abline(h=mean(fat$ankle[100:175]))  
> points(mean(fat$age[100:175]), mean(fat$ankle[100:175]),  
+ pch=16, col=rgb(.35,0,0))
```



Covariance

Covariance measures the difference between the two variables in the four regions. Suppose we have a dataset

$\{(x_i, y_i) : 1 \leq i \leq n\}$ then

$$\text{Cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- Data with strong linear relationship $(x_i - \bar{x})(y_i - \bar{y})$ will have the same sign. (i.e if data lies in first and third box or in second and fourth box).
- In such cases covariance will be large in absolute value.

Pearson Correlation Coefficient

Correlation is Covariance in standardised scale. Suppose we have a dataset $\{(x_i, y_i) : 1 \leq i \leq n\}$ then

$$\text{Cor}(x, y) = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{(x_i - \bar{x})}{S_x} \right) \left(\frac{(y_i - \bar{y})}{S_y} \right)$$

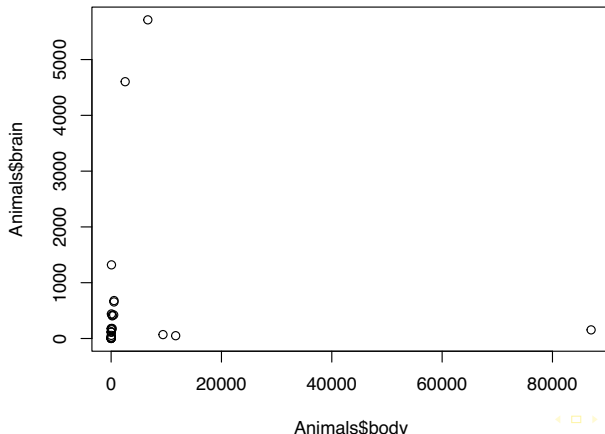
- $\text{Cor}(x, y)$ is between -1 and 1 .
- $\text{Cor}(x, y) \in \{1, -1\}$ indicates perfect linear relationship.
- $\text{Cor}(x, y) = 0$ indicates no linear relationship.

Paired Data: dataset in UsingR

```
> cor(fat$wrist, fat$neck)
[1] 0.7448264
> cor(fat$wrist, fat$height)
[1] 0.3220653
> cor(fat$age, fat$ankle)
[1] -0.1050581
```

Pearson Correlation Coefficient

```
> require(MASS)  
> plot(Animals$body,Animals$brain)
```



Package
- MASS

Dataset
- Animals

larger bodies



larger brains

Spearman Correlation Coefficient

```
> require(MASS)
> cor(Animals$body,Animals$brain)
[1] -0.005341163
```

- One way is to exclude the outliers.
- Another method is to transform the dataset by placing data in order and assigning a rank. Use `rank`.

```
> require(MASS)
> cor(rank(Animals$body), rank(Animals$brain))
[1] 0.7162994
```

or

```
> require(MASS)
> cor(Animals$body, Animals$brain, method="spearman")
[1] 0.7162994
```

Spearman Correlation Coefficient

$(X_i, Y_i) \quad i = 1, 2, \dots, n$ Data set

— Convert to rank of each data point

$(R(X_i), R(Y_i))$

$r_s \equiv$ Spearman Correlation Coefficient

$\text{Corr}(R(X), R(Y))$

Example:-

$x = 2, 3, 5, 7, 11$

Rank of x $R(x) = 1, 2, 3, 4, 5$

$y = 5, 5, 2, 7, 5$

Rank of y $R(y) = 3, 3, 1, 5, 3$

(ties have average rank)

$r_s = \text{Corr}(R(x), R(y))$

≈ 0.925

Spearman Correlation Coefficient

Suppose we have a dataset $\{(x_i, y_i) : 1 \leq i \leq n\}$ then first rank them to get $\{(r_{x_i}, r_{y_i}) : 1 \leq i \leq n\}$

$$\text{Spearman Correlation}(x, y) = \text{Cor}(r_x, r_y)$$

- measurement of relationship of monotonic data.
- not restricted to linear.

Chocolates and Noble Prizes

Correlation may NOT be Casual

Chocolate consumption and Nobel Prizes: A bizarre j—

<http://blogs.scientificamerican.com/the-curious-wave-...>

Chocolate consumption and Nobel Prizes: A bizarre j—

<http://blogs.scientificamerican.com/the-curious-wave-...>

nature PUBLISHING INDEX 2012 GLOBAL

Where does your institution rank?

Sign in / Register

Subscription Center

Subscribe to Print + Tablet

Subscribe to Print + e

Give a Gift

View the Latest Issue >

Subscribe News & Features Topics Blogs Videos & Podcasts Education Citizen Science SA Magazine SA Mind Products

About the SA Blog Network Choose a blog...

The Curious Wavefunction

Thoughts on chemistry and the history and philosophy of science

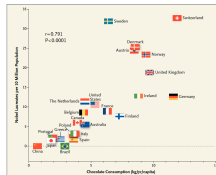
The Curious Wavefunction Home

Chocolate consumption and Nobel Prizes: A bizarre juxtaposition if there ever

WAVE

By Andrew J. Waite November 20, 2012 1

Share via Email Print



Correlation of chocolate consumption with Nobel Laureates (Image credit: New England Journal of Medicine)

What makes a Nobel Prize winner? There's several suggested factors: Perseverance? Good luck? Good mentors and students? Here's one possible

More from Scientific American

MIND > CHOCOLATES > DIGITAL >

THE ANIMAL LIVER'S DARK SIDE

Blog Network Highlights

The Thoughtful Animal

The Animal Liver's Dark Side

GIVE A GIFT - GET A GIFT!

GIVE A GIFT A GET OUR 2014 HUBBLE TELESCOPE CALENDAR FREE!

ORDER NOW!

Must Read Picks

Latest Picks

Peppino: Air Pollution enters from Beijing to Shanghai, at sea level again

Cholesterol: How Much Protein the Math Theorists behind the Fastest Flamingo Building

factor that I would have never imagined in my wildest dream: chocolate consumption. Chocolate consumption tracks well with the number of Nobel Laureates produced by a country.

At least that's what a paper published in the New England Journal of Medicine - one of the world's premier journals of medical research - claims. I have to say I found the study bizarre when I read it, and a few hours of strenuous, perplexed thought have done nothing to shake that feeling off. The study itself is amusing and rather brief and I think it makes for entertaining reading; what I am left contemplating is why this paper constitutes serious research and why it would have been published in a journal which over the years has presented some of the definitive medical findings of our time.

The paper starts by assuming - entirely reasonably - that winning a Nobel Prize must somehow be related to cognitive ability. It then goes on to describe a link between flavanols - organic molecules found among other foods in chocolate, green tea and red wine - and cognitive ability. Now I haven't read the literature on flavanols and cognitive ability, but I am sure that flavanols themselves couldn't possibly be responsible for improved cognitive effect, especially when they are part of a complex cocktail of dietary and environmental factors affecting brain function.

But let's say that's true; flavanols are indeed a strong indicator of cognitive function. From this idea the author basically jumps to the dubious and frankly bizarre question of whether chocolate consumption could possibly account for Nobel Prize winning ability. However, from a purely scientific standpoint the hypothesis is testable, so the author decides to simply plot the number of Nobel prize winners per 10 million people in different countries counted from 1900-2011 vs the chocolate consumption in those countries. The figures for chocolate consumption come from Cadkin and Chocoman and cover only four years, none before 2002. This fact itself makes any such comparison dubious to say the least; how can you compare two variables when they are sampled from such radically dissimilar sample spaces? And what about other compounds containing flavanols? We not also consider red wine or green tea?

In any case, a plot of chocolate consumption vs number of Nobel Prizes reveals a strong correlation of 0.79. Sweden is an anomaly (and the author thinks it could be a result of "paranoid bias" from the Nobel Committee); take it out and the correlation improves to 0.86. The graph in all its glory is illustrated above.

What does one make of this? Well, I have said before that if only three rules of scientific deduction were inscribed on the doors of every university and research organization in the world, one of them should be that "correlation does not mean causation". Conflating the two can lead you to believe, for instance, that *storks deliver babies*. Now the author recognizes this, but what I find absolutely baffling is that he makes no attempt to dissect other possible contributing factors. In fact at the end of the article he acknowledges the existence of such factors and then proceeds to dismiss them by saying that "differences in socioeconomic status from country to country and geographic and climatic factors may play some role, but they fall short of fully explaining the close correlation observed."

Observations:

Are Green Really Healthier? [Slide]

Observations:

Can We Avert the End of Elkharts?

Observations:

The ENIGMA Address on Agricultural Antibiotics in Overuse - and Utterly Ineffective

Follow Us:

See what new writing about Scientific American Editors

Free Newsletters

Get the best from Scientific American in your inbox

Email address

Tap into your MIND

GET BOTH Print + Tablet Editions of the Science of U.S. Energy: A Q&A with Secretary Ernest J. Moniz

Latest Headlines on ScientificAmerican.com

The Life of Dan Morgan: A Working Place for Joy and Intelligence

Energy Secretary Moniz to create US on trade mission (London)

The Science of U.S. Energy: A Q&A with Secretary Ernest J. Moniz

Hannovering Disruption to Move Forward with Blockchain (Prague)

The Animal Liver's Dark Side

Latest from

Children Make Highlight Lightbulb Best for Women in Research Centers

Let's Know, Not Let's Guess, We'll do it

Artificial Aliveness: Why Therapists Were Light-headed

Why Science Blogging Requires Story-Telling

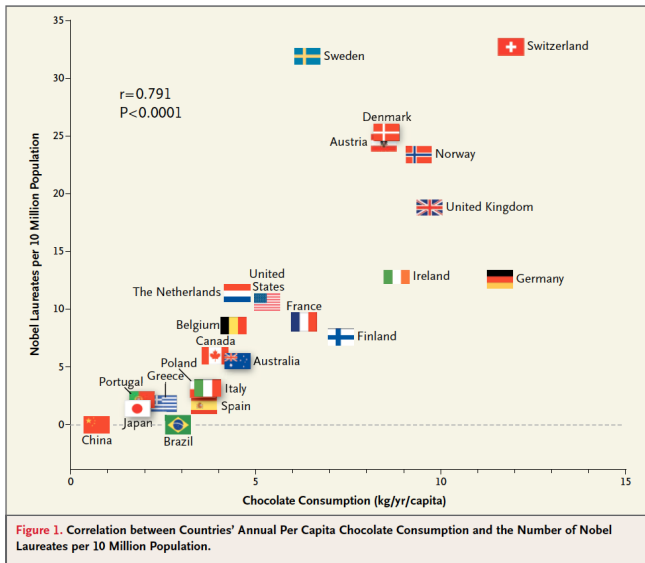
What is Science Blogging?

ADVERTISING

Scientific American In-Depth Reports

Comprehensive look at timely topics through with articles, podcasts, and interactive media.

Chocolates and Noble Prizes

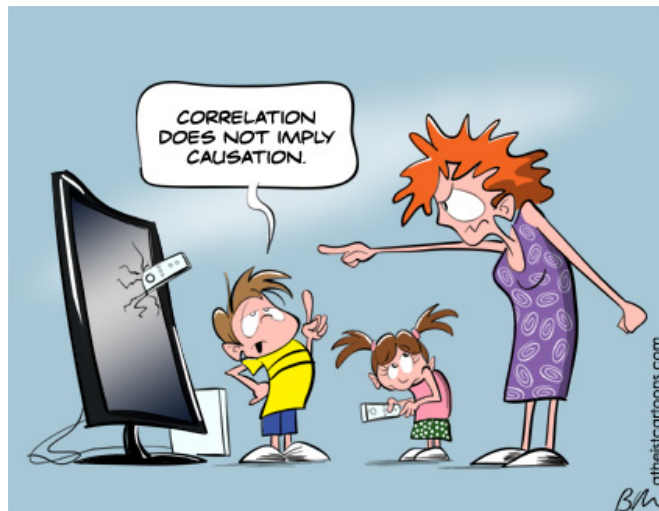


Chocolates and Noble Prizes

Noticed: Countries with more per capita chocolate consumption have more per capita Nobel laureates.

Conclude: Chocolate consumption cause better scientific research !

Chocolates and Noble Prizes

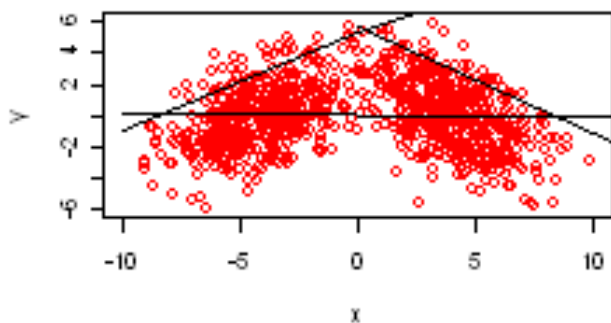


Chocolates and Noble Prizes

- Spurious: Facebook Users and Marks of users
- Causality: Smoking and lung cancer, Wine and heart risk.

Correlation

- Non-linear relationship
- 0 correlation



Correlation

- Pearson correlation coefficient is a measure of the linearity of the (possible) relationship between two variables X and Y .
- Even if correlation coefficient is high, it does not mean there is causal relationship between X and Y . Does not tell you cause and effect ?
- Care to be taken when used for predictive purposes.
- Causality: Domain Knowledge, design a good control experiment.

Basic Model [Simple linear Regression]

- Given data set $(X_i, Y_i) \quad i=1, \dots, n$
- Is there a relationship?
(linear)

Goal:- Model should provide an accurate low dimension summary.

Two parts :- - Define a family of models

Fitted model
"best"
among family
of models
- "NOT TRUTH"

- express a precise relationship
between Y and X

- Generate a fitted model

- closest model from the family
of models for the dataset

- "All models are wrong but some are useful"
- part of a larger text

$PV = RT$ - model for ideal gases
- rarely observed in nature

Simple linear Regression

- Given data set $(X_i, Y_i) \quad i=1, \dots, n$
- Is there a linear relationship?

$$"y_i = \beta_0 + \beta_1 x_i" \text{ for some } \beta_0, \beta_1 \in \mathbb{R}$$

Model:

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

Error:

$$E[\varepsilon] = 0$$

$$\text{Var}(\varepsilon) = \sigma^2$$

$$\varepsilon \perp X$$

$$\Rightarrow E[Y|X=x] = \beta_0 + \beta_1 x + 0$$

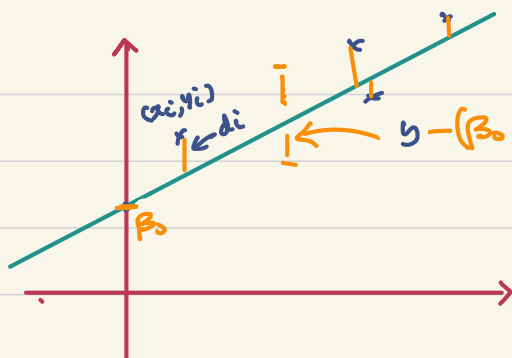
For data
 \Rightarrow
set

$$y_i = \beta_0 + \beta_1 x_i \quad (i=1, \dots, n)$$

Questions:

- ① Is the relationship linear?
- ② How to estimate β_0, β_1 ?
- ③ Can we provide confidence intervals for β_0, β_1 ?

$$\varepsilon \stackrel{d}{=} N(0, \sigma^2)$$



$$y = \beta_0 + \beta_1 x$$

$$y - (\beta_0 + \beta_1 x) \approx \text{Normal}(0, \sigma^2)$$

$$\beta_0 = E[y | x=0]$$

β_1 = slope of the line

$$\sigma^2 = \text{var}(y | x=x)$$

Method of Least Squares

$$\text{let } d_i = (y_i - \beta_0 - \beta_1 x_i)$$

- Find β_0, β_1 that minimize

$$\sum_{i=1}^n d_i^2 \equiv \text{Residual sum of squares}$$

- Calculus or last class method

to conclude

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}$$

Simple Linear Regression: Relationship in Bivariate Data

- Key: conditional mean of response variable given the predictor variable is a linear function.
- Model: For data points (x_i, y_i) with $1 \leq i \leq n$,

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i,$$

where ε_i assumed to be mean 0 and variance σ^2 Normal random variables.

- Observe only (x_i, y_i) for $1 \leq i \leq n$.

Simple Linear Regression: Relationship in Bivariate Data

- Find β_0, β_1 such that

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

is minimized.

- Can be solved: Calculus and Linear Algebra

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \text{correlation}(x, y) \frac{S_x^2}{S_y^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Observations:

- Slope of line is function of Correlation in standardised scale. 

Simple Linear Regression

[1] 0.7448264

$x := \text{wrist}$
 $y := \text{neck}$

$$y = \hat{\beta}_0 + \hat{\beta}_1 x$$

