

# Linear Models

Statistics with R  
Basel R Bootcamp



April 2019

# Linear Models

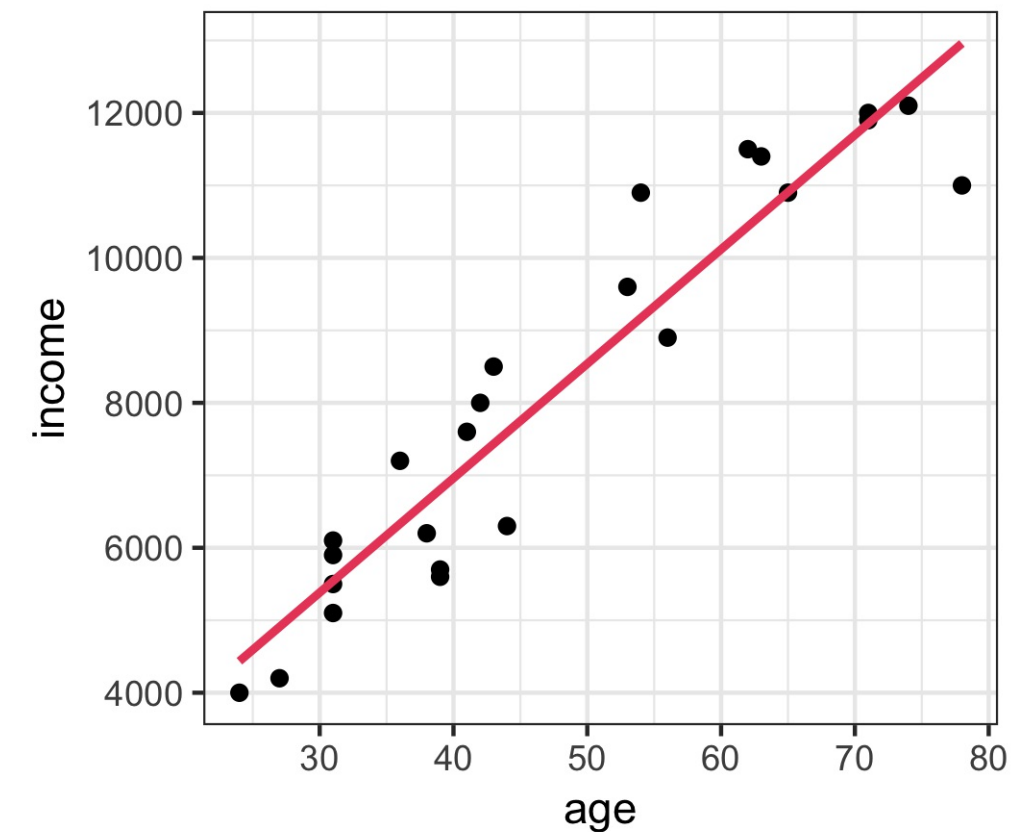
Linear models are, by far, **the most important models in all of statistics**.

Many statistical tests you may know of are special types of linear models.

Why are linear models so great?

- They are **easy to interpret**.
- They can **approximate non-linear** data well.
- They are **easy to calculate** and implement (just addition and multiplication!).
- They **just work**!

$\text{Large income} = 885 + \text{age} \times 149.3$



# What is a linear model?

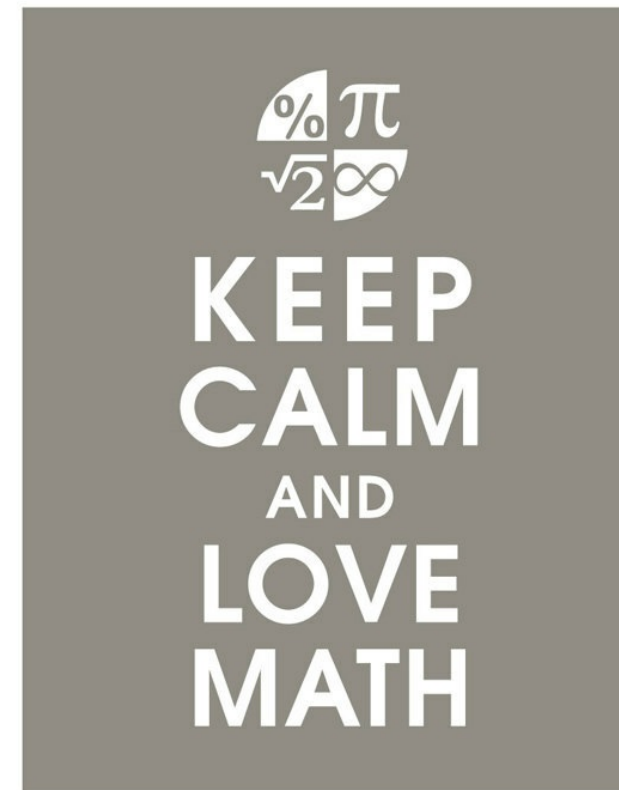
A linear model is **just addition and multiplication** and can be written in the following forms:

## Version 1

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

## Version 2

$$y = \beta_0 + \sum_{i=1}^n \beta_i x_i + \epsilon$$



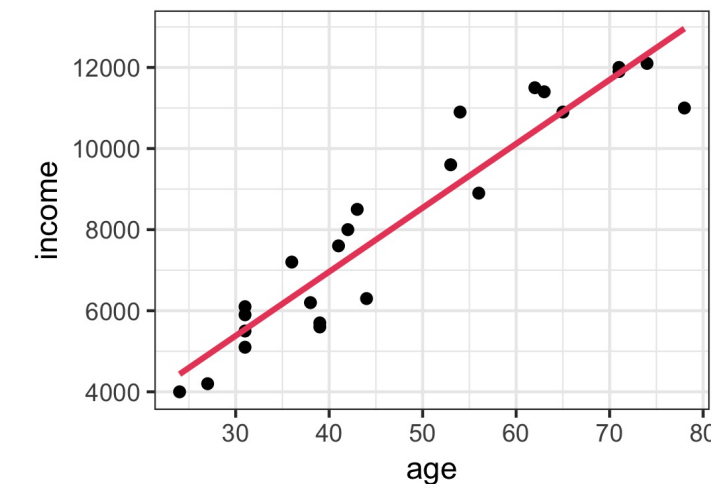
from [media.tumblr.com](https://www.tumblr.com/media.tumblr.com)

# Simple Linear Regression

**Definition:** Simple linear regression is a linear model with one predictor  $(x)$ , and where the error term  $(\epsilon)$  is Normally distributed.

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

Parameter	Description	In words
$\beta_0$	Intercept	When $x = 0$ , what is the predicted value for $y$ ?
$\beta_1$	Coefficient for $x$	For every increase of 1 in $x$ , how does $y$ change?



## Formula

$$\text{Large income} = 885 + 149.3 \times \text{age} + \epsilon$$

## Coefficients

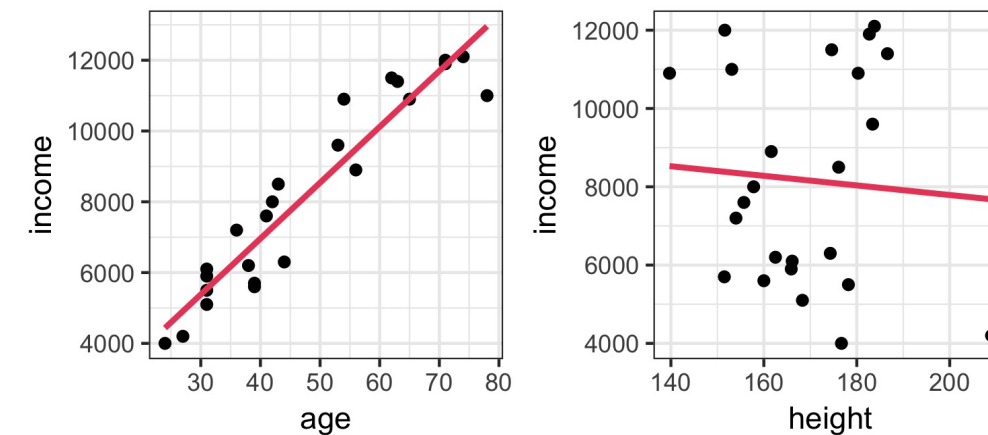
$$\text{Large } \beta_0 = 885, \beta_{\text{age}} = 149.3$$

# Multiple Linear Regression

**Definition:** Multiple linear regression is a linear model with many predictors  $(x_1, x_2, \dots, x_n)$ , and where the error term  $(\epsilon)$  is Normally distributed.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

Parameter	Description	In words
$\beta_0$	Intercept	When all x values are 0, what is the predicted value for y?
$\beta_1, \beta_2, \dots$	Coefficient for $x_1, x_2, \dots$	For every increase of 1 in coefficient for $x_1, x_2, \dots$ how does y change?



## Formula

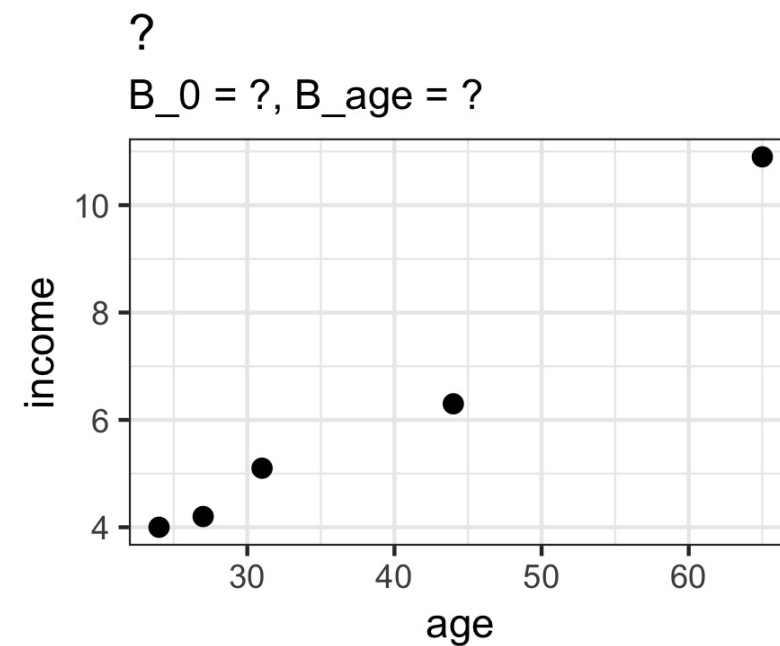
$$\text{income} = 1628 + 147 \times \text{age} - 4.1 \times \text{height} + \epsilon$$

## Coefficients

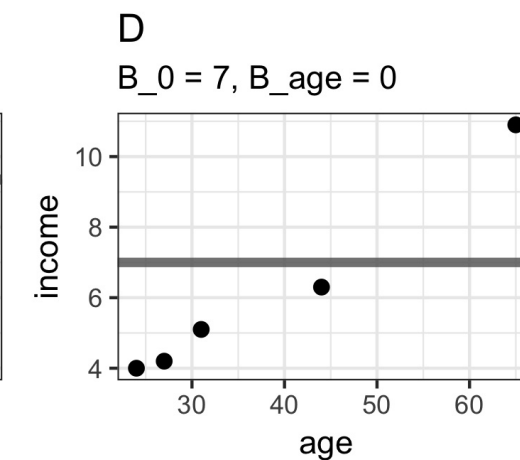
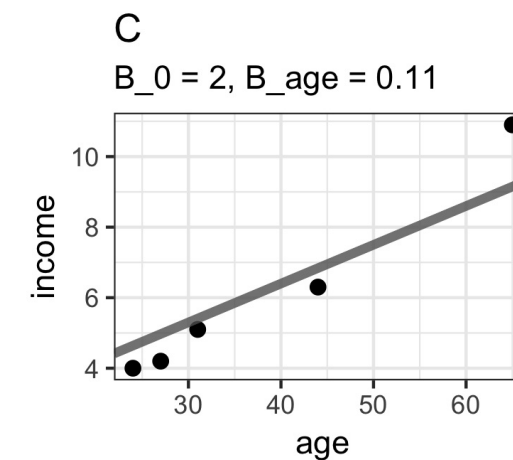
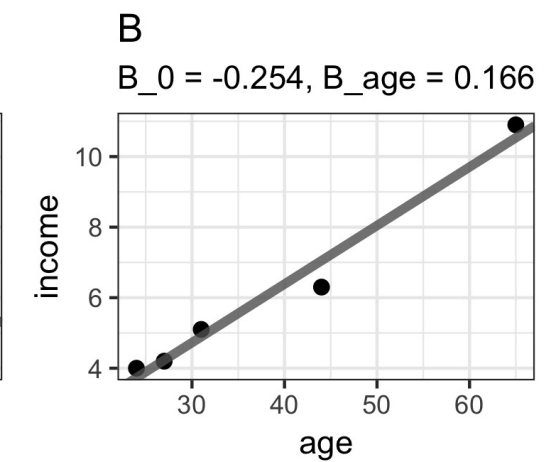
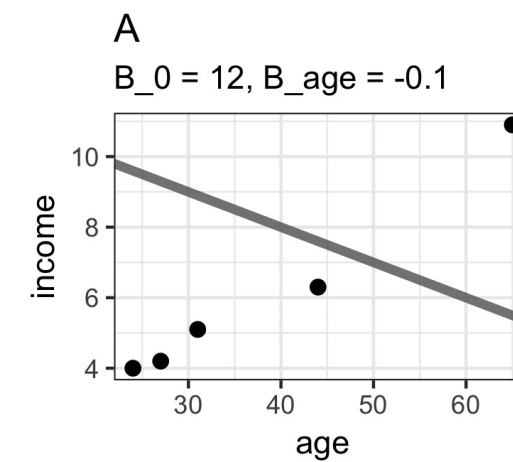
$$\beta_0 = 1628, \beta_{\text{age}} = 147, \beta_{\text{height}} = -4.1$$

# Estimating coefficients

How do we estimate the 'right' coefficients in a linear model? **Which of the four would be a good fit?**

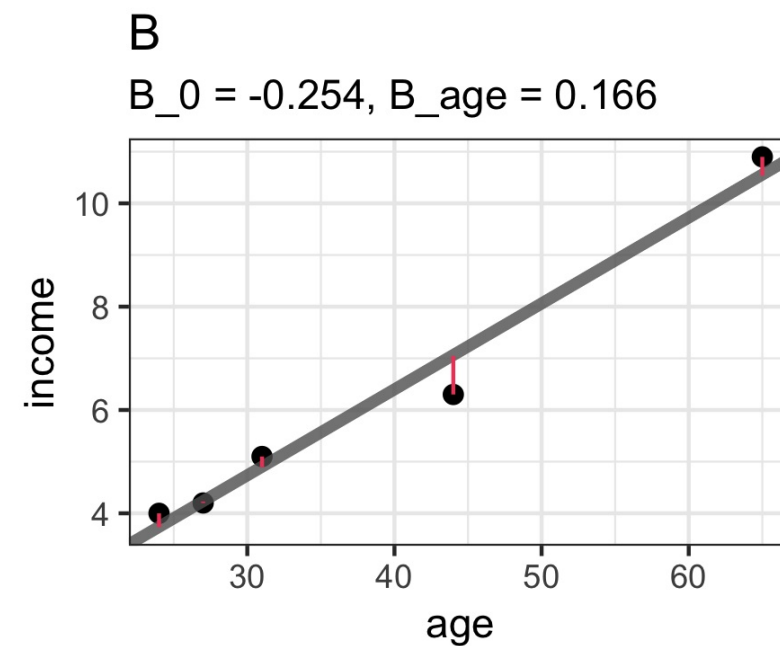


$$\text{Large income} = \beta_0 + \beta_1 \times \text{age} + \epsilon$$

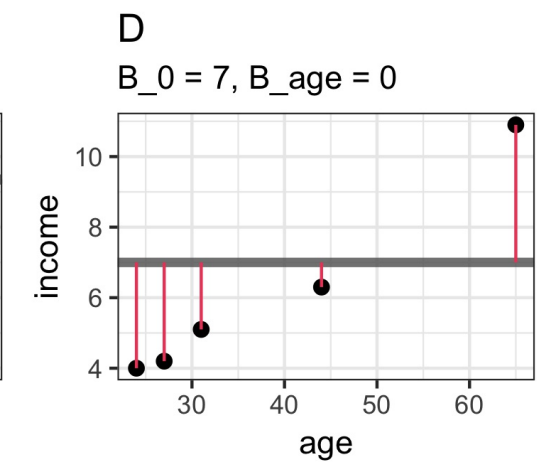
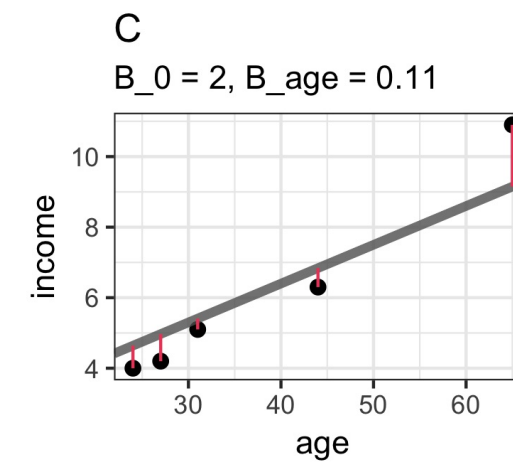
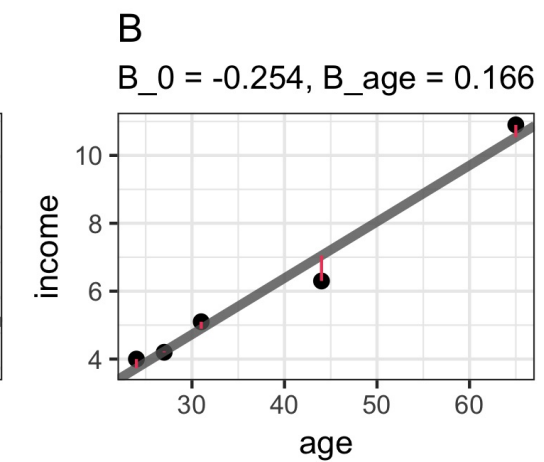
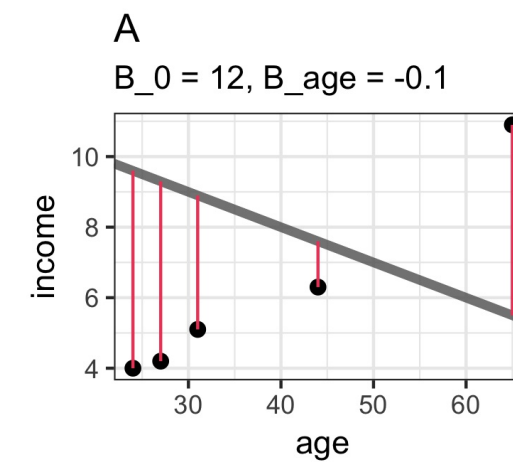


# Estimating coefficients

How do we estimate the 'right' coefficients in a linear model? **Which of the four would be a good fit?**



$\text{Large income} = -0.254 + 0.167 \times \text{age} + \epsilon$



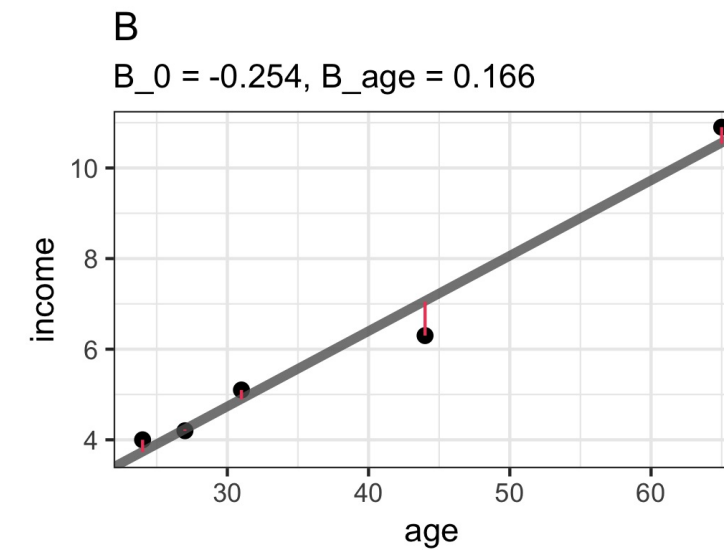
# Estimating coefficients

How do we estimate the 'right' coefficients in a linear model?

Find the values that minimise **Mean Squared Error** (MSE).

MSE: The average squared distance between the prediction and the data.

$$\text{Large MSE} = \frac{1}{N} \sum_{i=1}^n (y_i - \text{Prediction}_i)^2$$



id	age	income	Prediction	SE
1	24	4.0	3.730	0.0729
2	27	4.2	4.228	0.0008
3	31	5.1	4.892	0.0433
4	44	6.3	7.050	0.5625
5	65	10.9	10.536	0.1325
			MSE	0.16

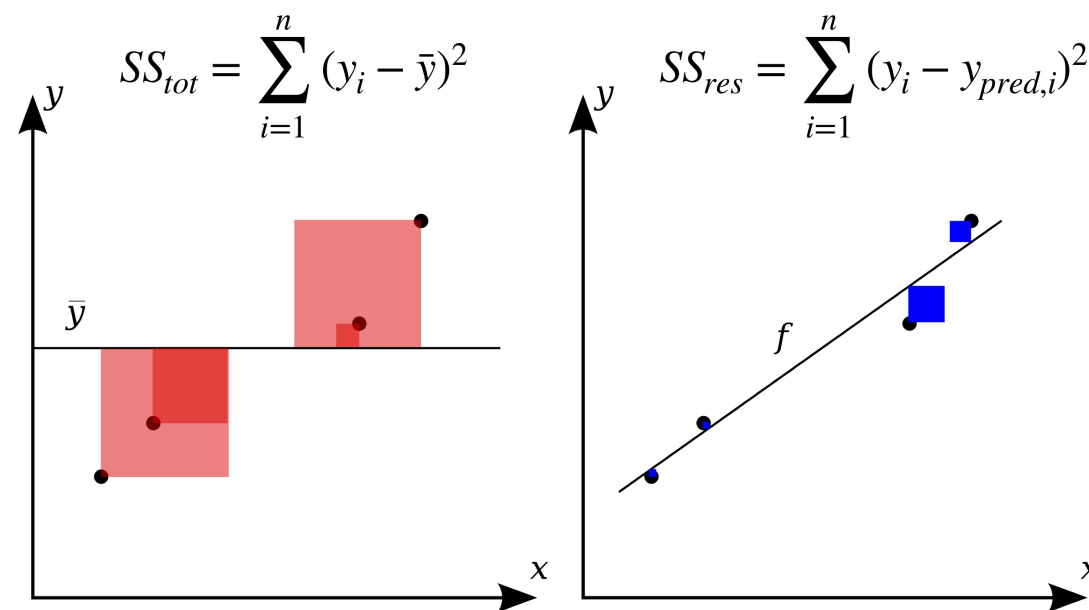


# R<sup>2</sup> (R-Squared)

R-Squared ( $R^2$ ) is the most common method of calculating the **overall performance** of a model.

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

$R^2$	Interpretation
0	Model explains no variance in y.
.5	Model explains half the variance in y.
1	Model explains all the variance in y.



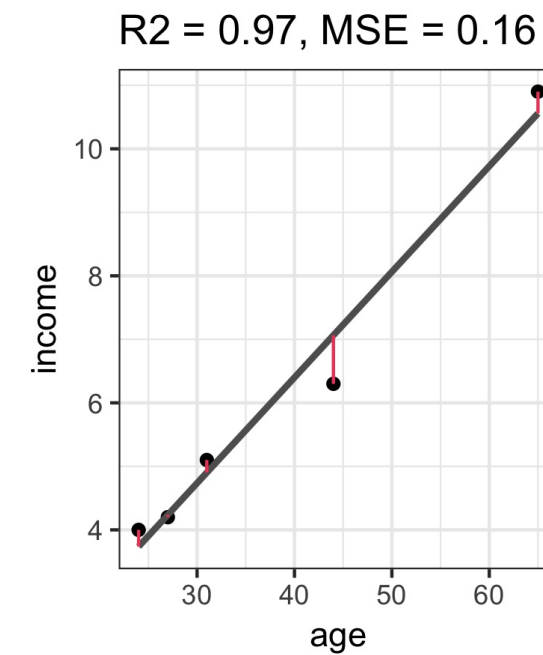
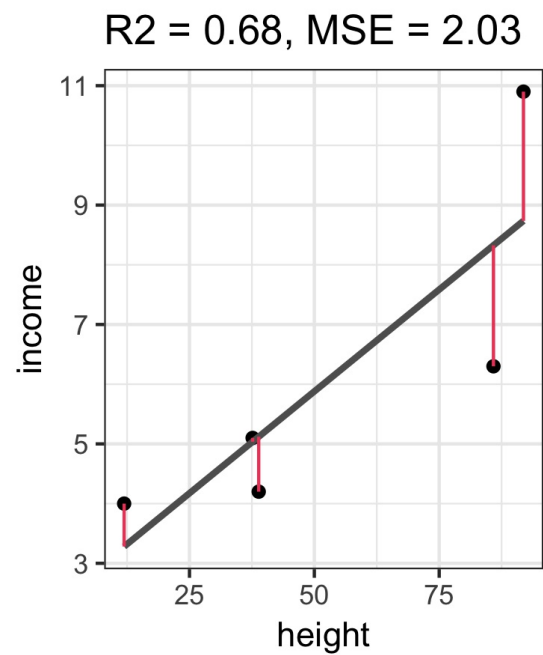
from Wikipedia

# R<sup>2</sup> (R-Squared)

R-Squared ( $R^2$ ) is the most common method of calculating the **overall performance** of a model.

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

R <sup>2</sup>	Interpretation
0	Model explains no variance in y.
.5	Model explains half the variance in y.
1	Model explains half the variance in y.



# How to fit and dissect linear models in R

## with `glm()`

# Key Functions

## Fitting

Function	Description
<code>glm(formula, data)</code>	Fit a linear model to data and calculate best coefficients

```
# Create a model income_glm

# Y = income
# X1 = age, X2 = children

income_glm <- glm(formula = income ~ age + children,
                  data = baselers)
```

## Evaluation

Function	Description
<code>coef(mod)</code>	Get coefficients from a model
<code>fitted(mod)</code>	Get fitted results.
<code>resid(mod)</code>	Get residuals (errors)

id	income	age	children
1	6300	44	2
2	10900	65	2

# Key Functions

## Fitting

Function	Description
<code>glm(formula, data)</code>	Fit a linear model to data and calculate best coefficients

## Evaluation

Function	Description
<code>coef(mod)</code>	Get coefficients from a model
<code>fitted(mod)</code>	Get fitted results.
<code>resid(mod)</code>	Get residuals (errors)

```
# Print income_glm
```

```
income_glm
```

```
##  
## Call:  glm(formula = income ~ age + children, data = baselers)  
##  
## Coefficients:  
## (Intercept)          age      children  
##      871.10       149.25         7.78  
##  
## Degrees of Freedom: 8509 Total (i.e. Null);  8507 Residual  
## (1490 observations deleted due to missingness)  
## Null Deviance:      6.33e+10  
## Residual Deviance: 1.29e+10    AIC: 145000
```

# Key Functions

## Fitting

Function	Description
<code>glm(formula, data)</code>	Fit a linear model to data and calculate best coefficients

```
# Show summary info
summary(income_glm)
```

## Evaluation

Function	Description
<code>coef(mod)</code>	Get coefficients from a model
<code>fitted(mod)</code>	Get fitted results.
<code>resid(mod)</code>	Get residuals (errors)

```
##
## Call:
## glm(formula = income ~ age + children, data = baselers)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -4250    -835         5      820    4779
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  871.101    44.959    19.4   <2e-16 ***
## age         149.249     0.818   182.5   <2e-16 ***
## children      7.777    12.861     0.6     0.55
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 1513138)
```

# Key Functions

## Fitting

Function	Description
<code>glm(formula, data)</code>	Fit a linear model to data and calculate best coefficients

## Evaluation

Function	Description
<code>coef(mod)</code>	Get coefficients from a model
<code>fitted(mod)</code>	Get fitted results.
<code>resid(mod)</code>	Get residuals (errors)

```
# Get fitted values (only first 5)
fitted(income_glm)
```

```
# Get residuals (only first 10)
resid(income_glm)
```

```
##      1      2      3      4      5
## 7454 10588 5513 4916 4461
```

```
##      1      2      3      4      5
## -1153.6 312.2 -413.4 -716.4 -460.9
```

id	income	age	children	fitted	resid
1	6300	44	2	7454	-1153.6
2	10900	65	2	10588	312.2

# Key Functions

## Fitting

Function	Description
<code>glm(formula, data)</code>	Fit a linear model to data and calculate best coefficients

```
library(rsq)
library(broom)

# Show R-squared from model
rsq(income_glm)

## [1] 0.7967
```

## Evaluation

Function	Description
<code>rsq(mod)</code>	Get R <sup>2</sup> value from model.
<code>tidy(mod)</code>	Get tidy results from a model.

```
# Show 'tidy' results from my model
tidy(income_glm)

## # A tibble: 3 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)    871.      45.0     19.4 7.11e-82
## 2 age           149.      0.818    183.  0.
## 3 children       7.78     12.9     0.605 5.45e- 1
```



# Practical