



A scatterplot of weight (Y) versus height (X) for 247 physically active

```
> library(openintro)
```

$$y_i = \beta_0 + \beta_1 x_i + e_i$$

- Remark:** We capitalize Y_i in the equation to emphasize that it is a random variable



Conditional Mean and Variance Functions

- ▶ This expected value of Y when X takes a specific value x .

$$E(Y|X = x) = \beta_0 + \beta_1 x$$

For SLR, the conditional mean is modeled as a straight line.

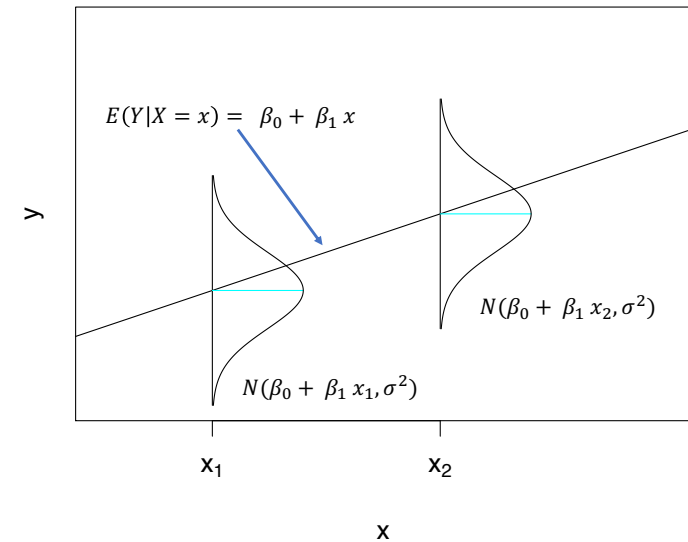
- ▶ The variance of Y when X takes a specific value x .

$$\text{Var}(Y|X = x) = \sigma^2$$

An assumption for SLR is that the variance is the same for every value of x .

- ▶ The conditional distribution of Y when X takes a specific value x .

$$Y|X \sim N(\beta_0 + \beta_1 x, \sigma^2)$$



Fitted Values and Residuals

The line that we estimate, or fit to the data in the scatterplot, is written as

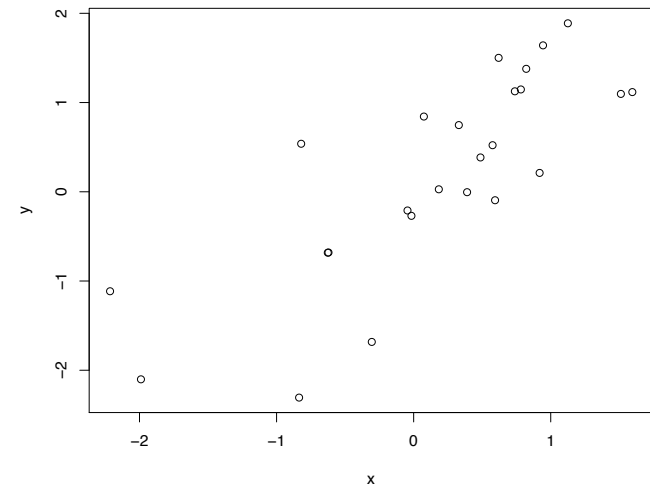
$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

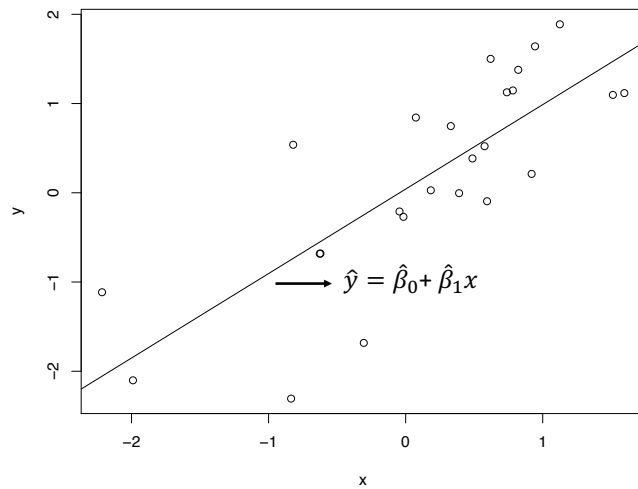
The fitted (or predicted) value for the i^{th} observation (x_i, y_i) :

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

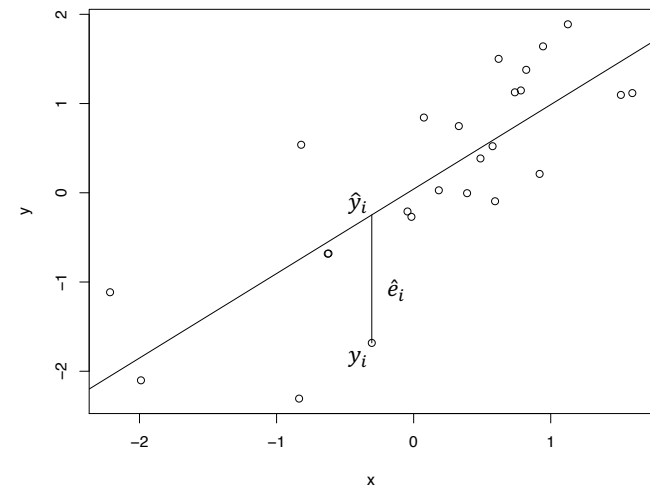
The **residual** for the i^{th} observation is the difference between the observed value (y_i) and the predicted value (\hat{y}_i):

$$\hat{e}_i = y_i - \hat{y}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$$



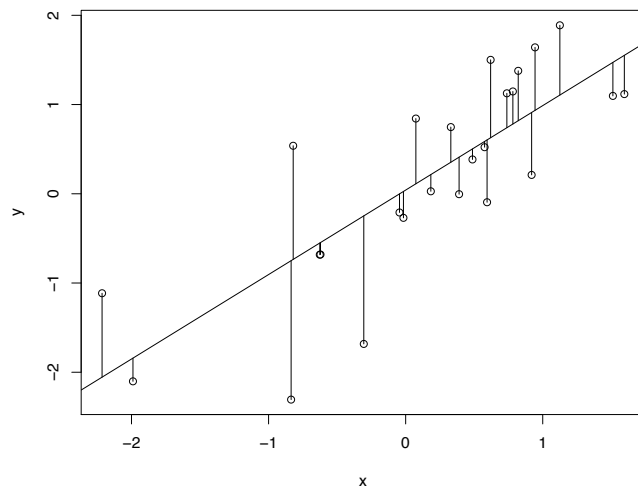


Navigation icons: back, forward, search, etc.



Navigation icons: back, forward, search, etc.

Sum of Squared Residuals



Navigation icons: back, forward, search, etc.

- Intuitively, a line that fits the data well has small residuals.
- The **least squares line** minimizes the **sum of squared residuals**:

$$RSS = \sum_{i=1}^n \hat{e}_i^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

- That is, out of all possible lines we could draw on the scatterplot, the least squares line is the “best fit” since it has the smallest sum of squared residuals.

Navigation icons: back, forward, search, etc.

Least Squares Estimation

Formally, the estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ of the intercept and slope are found by using calculus to minimize the sum of squared residuals:

$$RSS = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

To minimize set the partial derivatives equal to zero:

$$\frac{\partial RSS}{\partial \hat{\beta}_0} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$\frac{\partial RSS}{\partial \hat{\beta}_1} = -2 \sum_{i=1}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

Navigation icons

Interpretation

- **Slope:** an increase in the explanatory variable (x) by one unit is associated with a change of $\hat{\beta}_1$ in the predicted response (\hat{y}).
- **Intercept:** the prediction for the response variable (\hat{y}) when the value for the explanatory variable is zero ($x = 0$). It may not make sense to try to interpret the intercept depending on the application.

Navigation icons

Least Squares Estimation

Solving these two equations gives the **least squares estimates** of the intercept and slope:

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

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

Note that the equation for the intercept guarantees the least squares line passes through (\bar{x}, \bar{y}) .

Navigation icons

Estimating σ^2

Estimate of $\text{Var}(e_i) = \sigma^2$:

$$\hat{\sigma}^2 = \frac{RSS}{n-2} = \frac{1}{n-2} \sum_{i=1}^n \hat{e}_i^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Remarks:

- $\sum_{i=1}^n \hat{e}_i = 0$
- $\hat{\sigma} = \sqrt{RSS/(n-2)}$ called the **residual standard error**
- The divisor is $n-2$ since two parameters β_0 and β_1 were estimated
- It can be shown that $\hat{\sigma}^2$ is an unbiased estimate of σ^2 , i.e., $E(\hat{\sigma}^2) = \sigma^2$

Navigation icons

Example

```
> lm1 <- lm(wgt ~ hgt, data=bdims_males)
> summary(lm1)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-60.95336	14.05436	-4.337	2.11e-05 ***
hgt	0.78257	0.07901	9.905	< 2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8.902 on 245 degrees of freedom

Multiple R-squared: 0.2859, Adjusted R-squared: 0.283

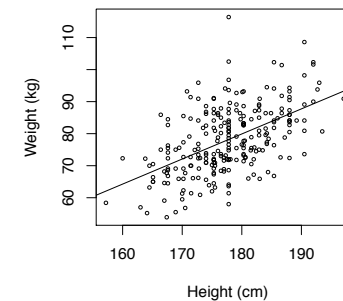
F-statistic: 98.11 on 1 and 245 DF, p-value: < 2.2e-16

Navigation icons

Example

A scatterplot of weight (Y) versus height (X) for 247 physically active men with least squares line superimposed.

```
> plot(wgt ~ hgt, data=bdims_males,
       xlab = 'Height (cm)', ylab = 'Weight (kg)', cex=0.5)
> abline(lm1)
```



Navigation icons

Example

- Equation of the least square regression line:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x = -60.95 + 0.78x$$

or we can write $\text{weight} = -60.95 + 0.78\text{height}$

- Interpreting the slope: For men, an increase in height by 1 cm is associated with an increase in weight by 0.78 kg.
- Interpreting the intercept: The predicted weight for a man that is 0 cm tall is -60.95 kg. Note that it does not make sense to interpret the intercept in this context. The prediction is an extrapolation (heights for men in this data set range between 157.2 to 198.1 cm).

Navigation icons

Example

We can calculate the residual standard error manually in R. This value is also given in the regression summary.

```
> n <- nrow(bdims_males)
> sqrt(sum(resid(lm1)^2) / (n-2))
[1] 8.901667
```

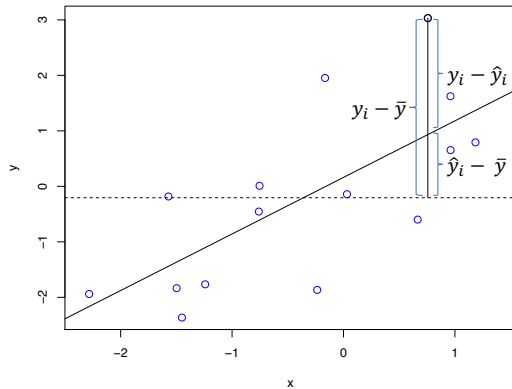
Note that the empirical sum of the residuals is approximately 0:

```
> sum(resid(lm1))
[1] -1.273981e-13
```

Navigation icons

Partitioning Variability

Graphical description that $y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$



Navigation icons: back, forward, search, etc.

Partitioning Variability

Remarkably, it can be shown that

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

$$SST = SS_{\text{reg}} + RSS$$

- ▶ SST is the total sum of squares (total variability in the response variable)
- ▶ SSreg is the regression sum of squares (variability in the response explained by the model)
- ▶ RSS is the residual sum of squares (unexplained variability)

Navigation icons: back, forward, search, etc.

Coefficient of Determination

The **coefficient of determination** (R^2) is a measure of how well the linear regression model fits the data.

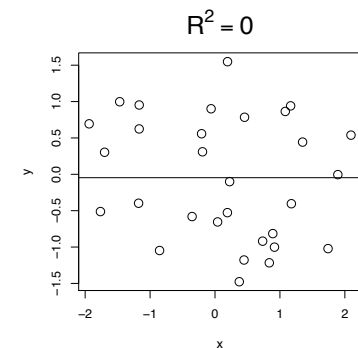
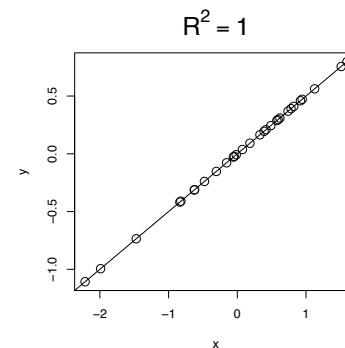
$$R^2 = \frac{SS_{\text{reg}}}{SST} = 1 - \frac{RSS}{SST}$$

- ▶ R^2 can be interpreted as the proportion of variability in the response variable Y that is explained by X (i.e., the regression model).
- ▶ $0 \leq R^2 \leq 1$; the closer R^2 is to 1, the better the linear regression model fits the data.
- ▶ For the example, $R^2 = 0.286$ (see summary output), meaning that for men, 28.6% of the variability in weight can be explained by height.

Navigation icons: back, forward, search, etc.

Limiting cases:

- ▶ $R^2 = 1$ when all points fall on the regression line (RSS=0)
- ▶ $R^2 = 0$ when $\hat{y}_i = \bar{y}$, which implies RSS=SST.



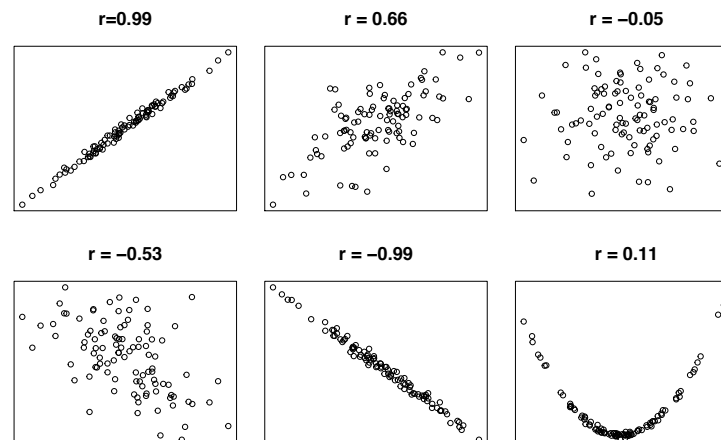
Navigation icons: back, forward, search, etc.

Correlation Coefficient (review)

The **correlation coefficient**, denoted by r , is a number between -1 and 1 that describes the strength of the linear association between two variables.

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

- ▶ \bar{x} and \bar{y} are the sample means
- ▶ s_x and s_y are the sample standard deviations

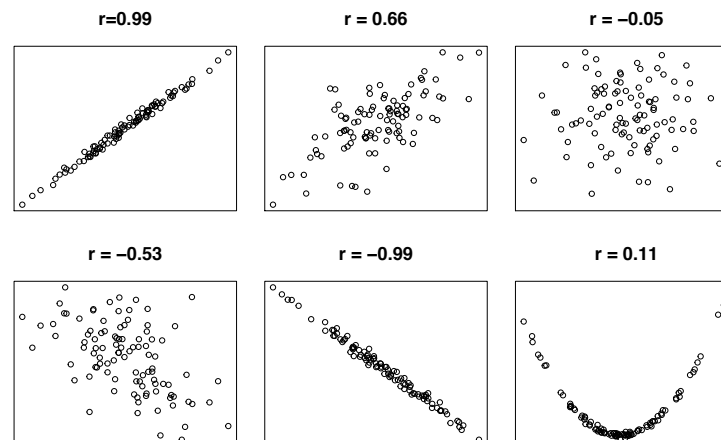


Correlation Coefficient

- ▶ $r \approx 1$ when there is a strong positive linear association between the variables.
- ▶ $r \approx -1$ when there is a strong negative linear association between the variables.
- ▶ $r \approx 0$ when there is no association between the variables (i.e., independent).
- ▶ The correlation coefficient is only useful for evaluating the linear association between two variables. It is not a useful measure for nonlinear relationships.



Correlation Coefficient



Correlation Coefficient

- ▶ R^2 can also be computed as the correlation coefficient squared.

```
> cor(bdims_males$wgt, bdims_males$hgt)^2
```

```
[1] 0.2859487
```
- ▶ The least squares estimate of the slope can be written in terms of the correlation coefficient:

$$\hat{\beta}_1 = r \frac{s_y}{s_x}$$

