

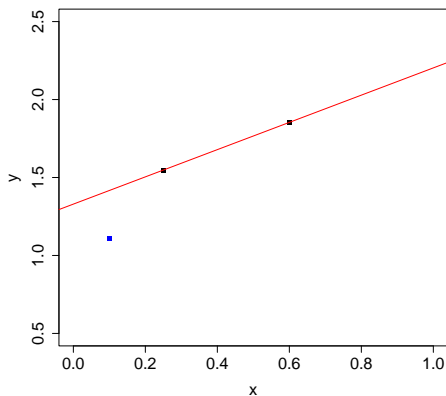
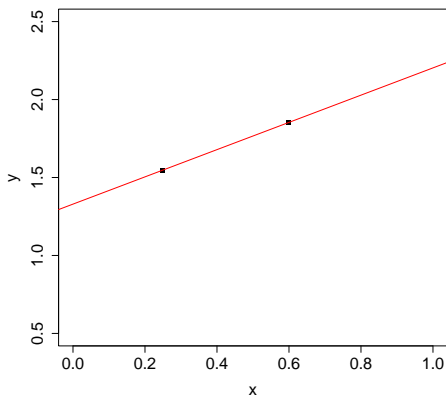
Regression Analysis

Chapter 1 and 2: Simple Linear Regression

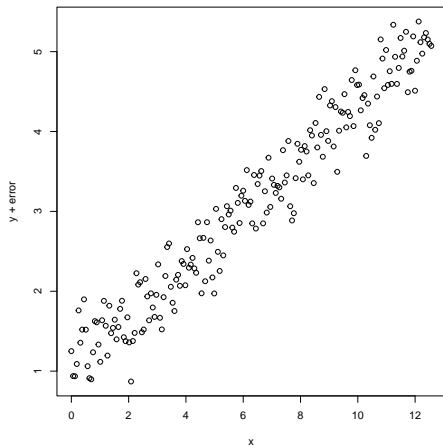
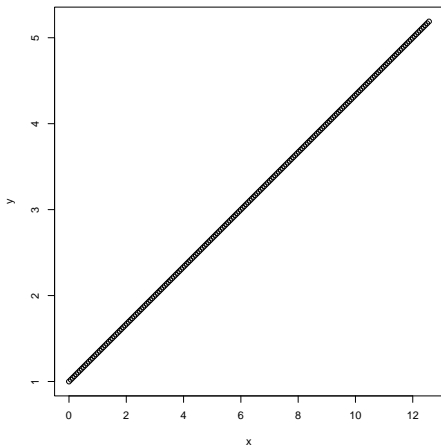
Shaobo Jin

Department of Mathematics

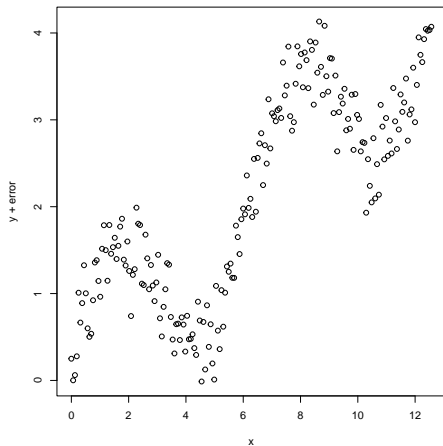
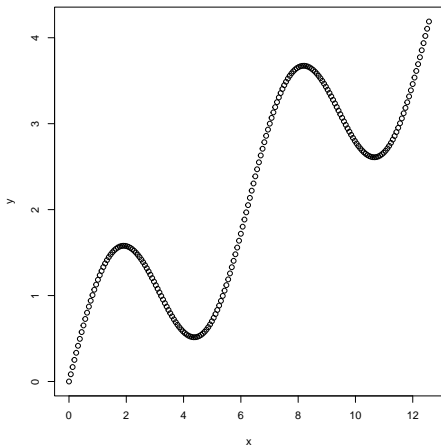
Two Points Determine A Line



Draw A Line/Curve



Draw A Line/Curve



Goal of Regression

The goal of regression is

- ① to summarize observed data as simply, usefully, and elegantly as possible,
- ② to make predictions.

Mean Function

The conditional expectation of the response variable Y when the [covariate/predictor/regressor](#) is fixed at $\mathbf{X} = \mathbf{x}$ is

$$E(Y \mid \mathbf{X} = \mathbf{x}) = f(\mathbf{x}),$$

which is called a [mean function](#). For example

- $f(x) = \beta_0 + \beta_1 x$ is linear in x , where β_0 (intercept) and β_1 (slope) are the parameters.
- $f(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$ is linear in \mathbf{x} and also linear in β 's, where β_0 is the intercept and β_1 and β_2 are the slopes.
- $f(x) = \beta_0 + \beta_1 x + \beta_2 x^2$ is not linear in x , but still linear in the parameters.
- $f(x) = \beta_0 + \beta_1 / [1 - \exp(\beta_2 x)]$ is not linear.

Variance Function

The conditional variance of the response variable Y when the predictor/regressor is fixed at $\mathbf{X} = \mathbf{x}$ is

$$\text{Var}(Y \mid \mathbf{X} = \mathbf{x}),$$

which is called a **variance function**.

- A frequent assumption is that $\text{Var}(Y \mid \mathbf{X} = \mathbf{x}) = \sigma^2$ (**homoscedasticity**).
- If $\text{Var}(Y \mid \mathbf{X} = \mathbf{x})$ depends on the value of $\mathbf{X} = \mathbf{x}$, then it is called **heteroscedasticity**.

Simple Linear Regression

We often write our model as

$$y = E(Y | X = x) + e,$$

where e is the random error term. One assumptions that we often make is

$$E(e | X = x) = 0.$$

The [simple linear regression](#) model consists of

$$\begin{aligned} E(Y | X = x) &= \beta_0 + \beta_1 x, \\ \text{Var}(Y | X = x) &= \sigma^2, \end{aligned}$$

where β_0 is the [intercept](#) and β_1 is the [slope](#).

Estimation: Assumptions

The parameters are β_0 , β_1 and possibly σ^2 . Suppose that we observe a data set (x_i, y_i) , $i = 1, 2, \dots, n$.

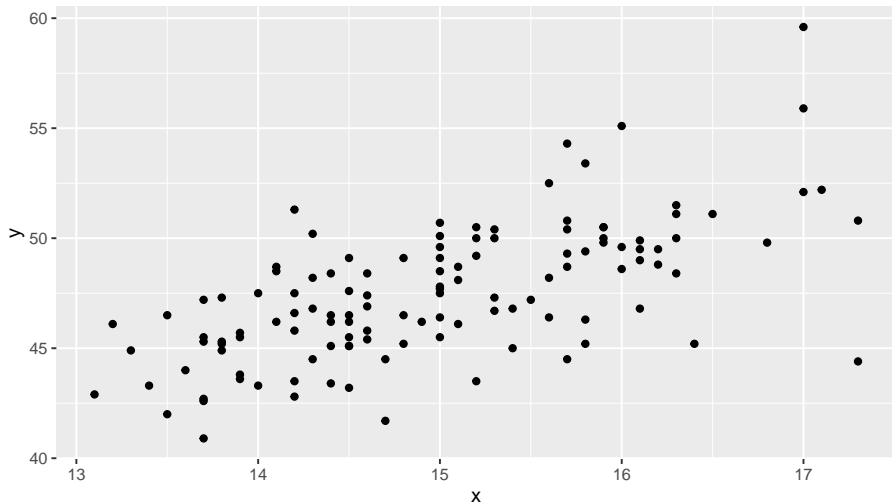
- We often write the simple linear regression model as

$$y_i = \beta_0 + \beta_1 x_i + e_i, \quad i = 1, 2, \dots, n.$$

- We often assume data are independent, i.e., given all (X_1, \dots, X_n) , e_i is independent of e_j for all $i \neq j$.

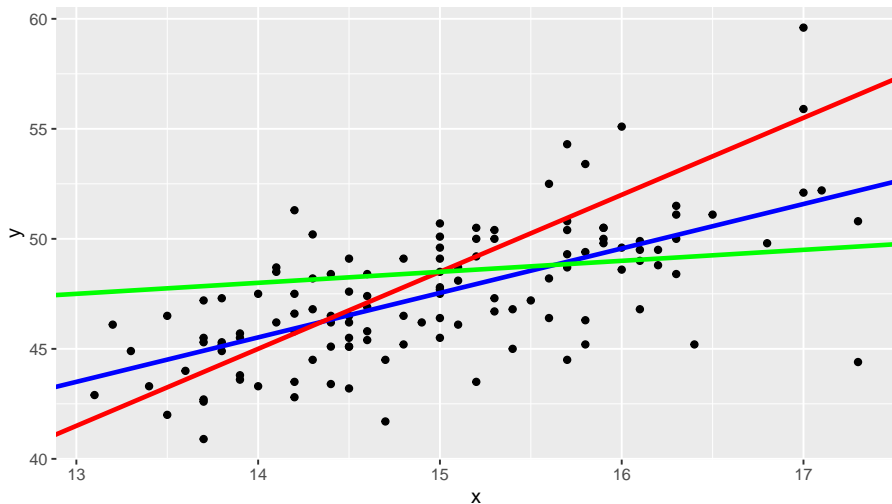
We need to estimate the unknown parameters using our data.

Scatter Plot



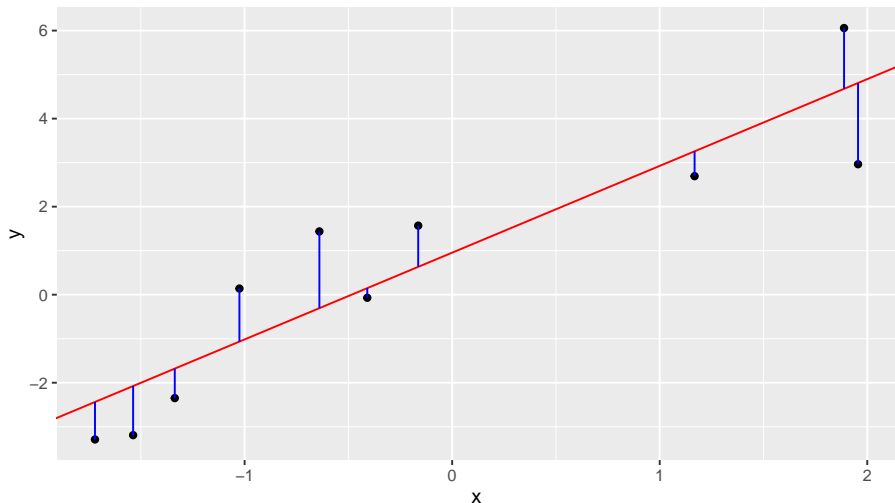
Scatter Plot

Which linear function describes the data the best?



Distance to Conditional Mean

For every straight line, we can compute the vertical distance of our data point to the line $y_i - (\beta_0 + \beta_1 x_i)$:



Ordinary Least Squares

The method of **ordinary least squares** (OLS) minimizes the sum of squares of $y_i - (\beta_0 + \beta_1 x_i)$. That is, the **OLS estimator** of β_0 and β_1 are the values such that

$$\text{RSS}(\beta_0, \beta_1) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$$

is minimized, where RSS stands for **residual sum of squares**. The minimizer is given by

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x},\end{aligned}$$

where $\bar{y} = \sum_{i=1}^n y_i / n$ and $\bar{x} = \sum_{i=1}^n x_i / n$.

Property of OLS Estimator: Linear Combination

In fact,

$$\begin{aligned}\hat{\beta}_1 &= \sum_{i=1}^n \left[\frac{(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] y_i, \\ \hat{\beta}_0 &= \sum_{i=1}^n \left[\frac{1}{n} - \frac{(x_i - \bar{x}) \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] y_i.\end{aligned}$$

Both are linear combinations of y_1, \dots, y_n .

Property of OLS Estimator: Unbiased

Under the assumption that $E(Y | X = x) = \beta_0 + \beta_1 x$ is correctly specified, the OLS estimators are unbiased as

$$E(\hat{\beta}_0 | \mathbf{x}) = \beta_0, \quad E(\hat{\beta}_1 | \mathbf{x}) = \beta_1,$$

where $\mathbf{x} = [x_1 \ x_2 \ \cdots \ x_n]^T$.

Property of OLS Estimator: Covariance

Under the assumptions that

① data are independent conditional on \mathbf{x} ,

② $\text{Var}(Y | X = x) = \sigma^2$,

the OLS estimators satisfy

$$\text{Var}(\hat{\beta}_0 | \mathbf{X}) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right],$$

$$\text{Var}(\hat{\beta}_1 | \mathbf{X}) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2},$$

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}_1 | \mathbf{X}) = -\frac{\sigma^2 \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Prediction/Fitted Value and Residual

Once β_0 and β_1 are estimated, the fitted/estimated regression line is

$$\hat{E}(Y \mid X = x) = \hat{\beta}_0 + \hat{\beta}_1 x.$$

The residual is

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

If we have a new observation with x_0 as covariate value, then the predicted value is

$$\hat{E}(Y \mid X = x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0.$$

Estimating σ^2

The estimate of σ^2 is

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{e}_i^2}{n-2}.$$

We call $\hat{\sigma}$ the [standard error of regression](#).

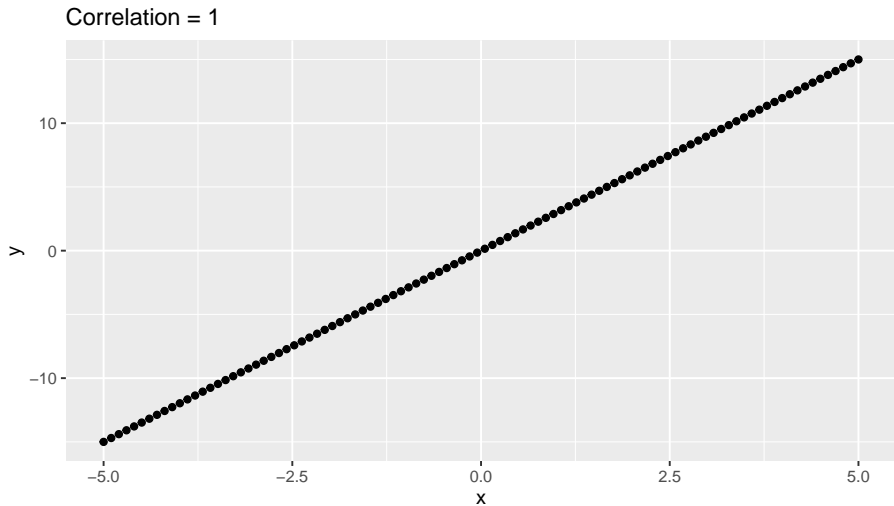
Correlation and Linear Regression

Note that

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\&= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}.\end{aligned}$$

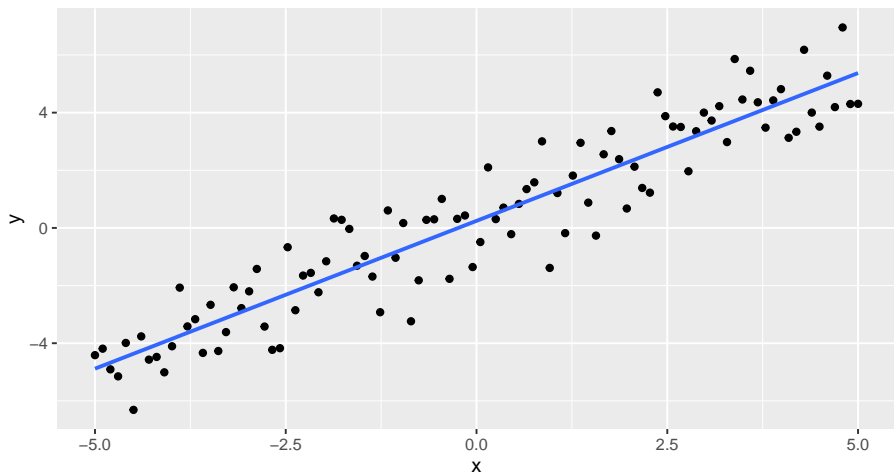
If the Pearson correlation is zero, then we don't have any linear trend. That is, Pearson correlation measure the degree of the linear relation between y and x .

Correlation and Linear Regression

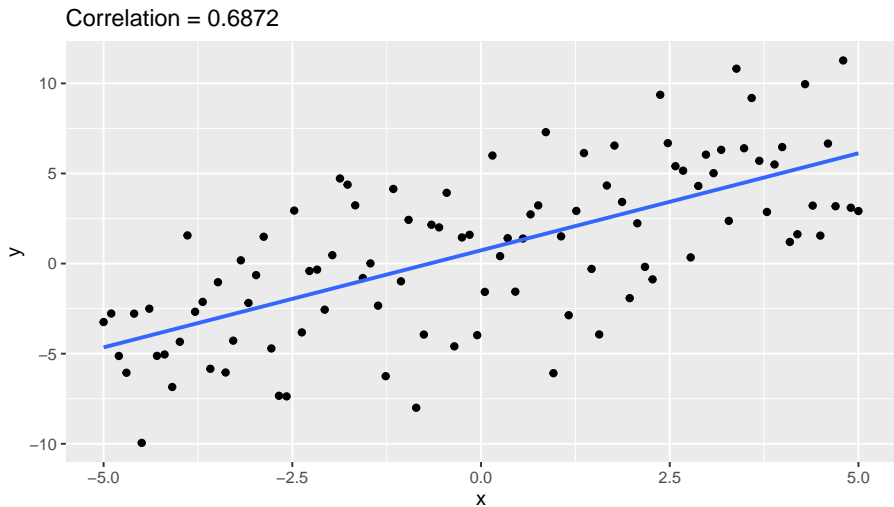


Correlation and Linear Regression

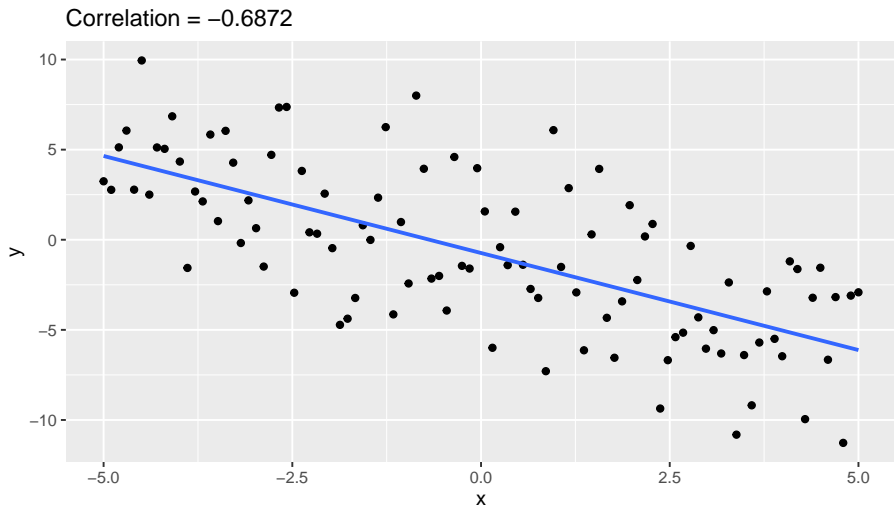
Correlation = 0.9379



Correlation and Linear Regression

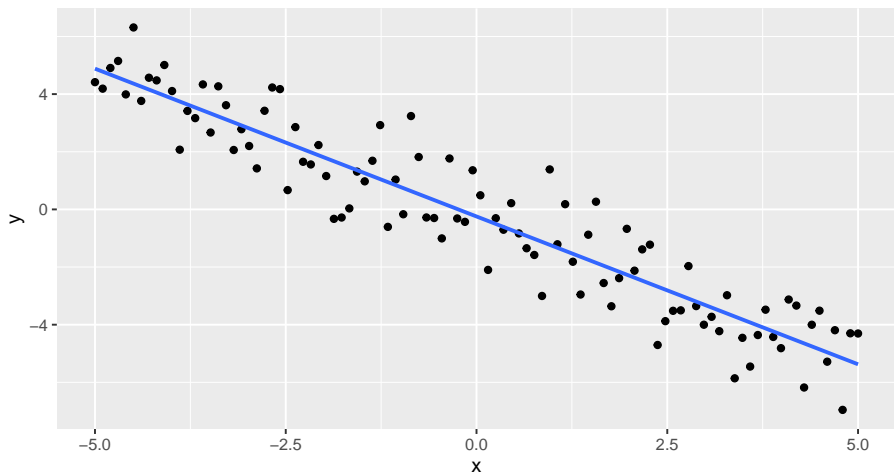


Correlation and Linear Regression

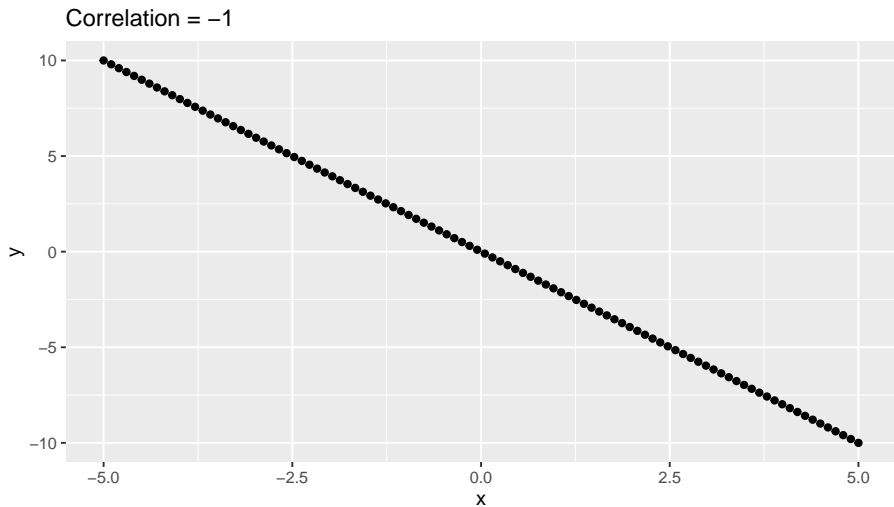


Correlation and Linear Regression

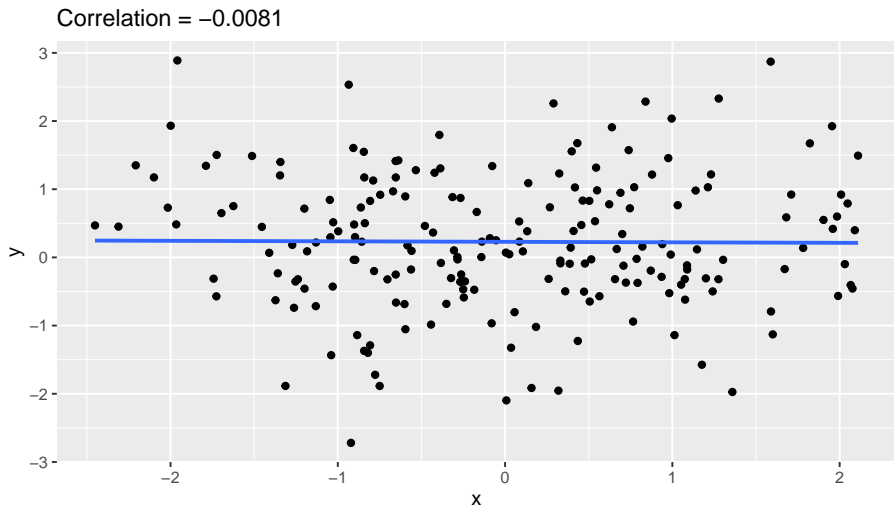
Correlation = -0.9379



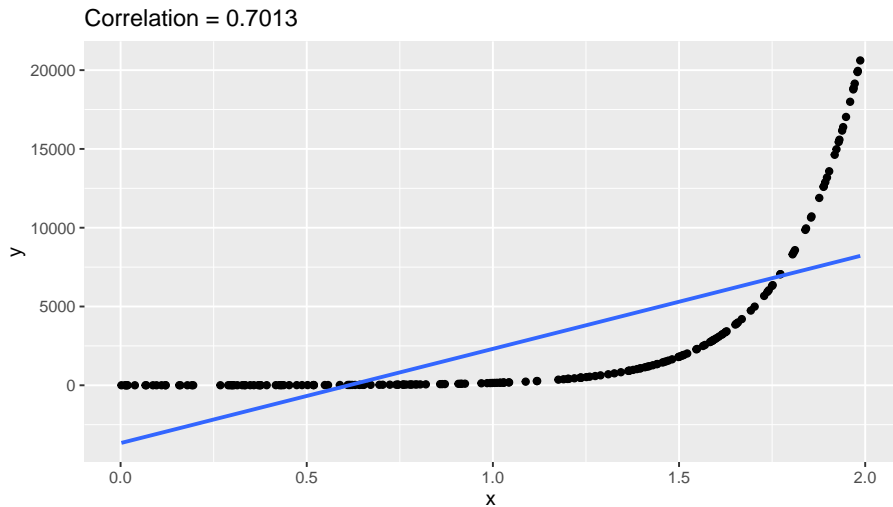
Correlation and Linear Regression



No Linear Relation



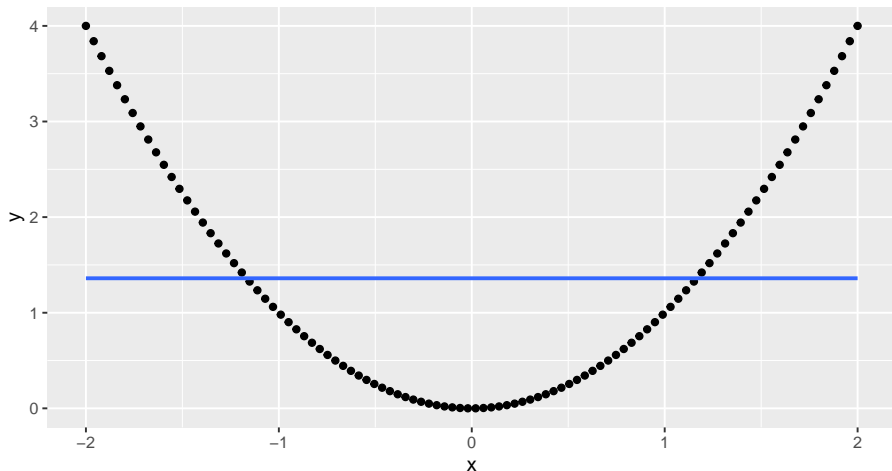
Nonlinear Relation



Problem with Pearson Correlation

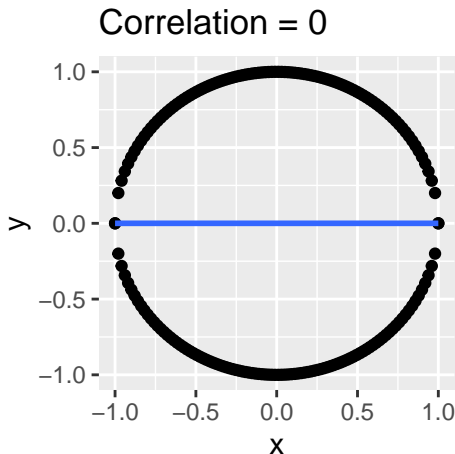
Pearson correlation only measure the degree of the linear relation between y and x .

Correlation = 0



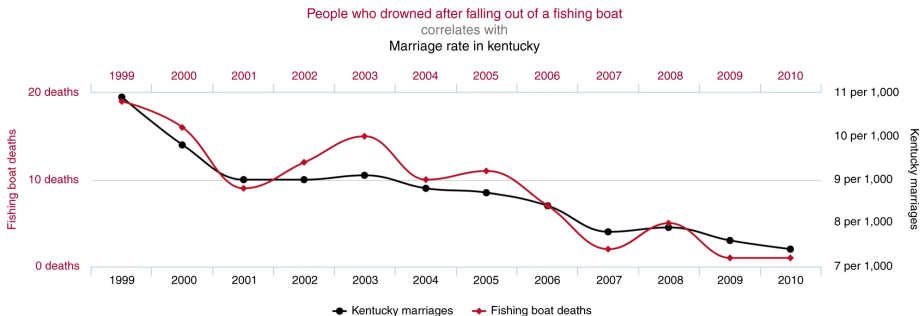
Problem with Pearson Correlation

Pearson correlation only measure the degree of the linear relation between y and x .



Correlation and Causation

Keep in mind that the regression model only describes correlation (association), not causation.



Normally Distributed e

We haven't assume the distribution of e yet. It is also common to assume

$$e \mid X = x \sim N(0, \sigma^2).$$

Under the independence assumption, the **log-likelihood function** of β_0 , β_1 , and σ^2 is

$$\begin{aligned} \ell(\beta_0, \beta_1, \sigma^2) = \sum_{i=1}^n \left\{ -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\sigma^2) \right. \\ \left. - \frac{1}{2\sigma^2} [y_i - (\beta_0 + \beta_1 x_i)]^2 \right\}. \end{aligned}$$

Maximum Likelihood Estimator

The **maximum likelihood estimator (MLE)** is given by

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x}, \\ \hat{\sigma}_n^2 &= \frac{1}{n} \sum_{i=1}^n \left[y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right]^2 = \frac{1}{n} \sum_{i=1}^n \hat{e}_i^2.\end{aligned}$$

The MLE of β_0 and β_1 are the same as their OLS estimator! Hence, they are still unbiased.

Random Vector and Random Matrix

A random vector/matrix is a vector/matrix of random variables. For example,

$$\mathbf{x} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix},$$
$$\mathbf{X}_{n \times p} = \begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix}.$$

Mean and Covariance Matrix

- ① The **expected value/mean** of a random vector/matrix is the vector/matrix consisting of the expected values of each of its elements:

$$E(\mathbf{X}_{n \times p}) = \begin{bmatrix} E(X_{11}) & \cdots & E(X_{1p}) \\ \vdots & \ddots & \vdots \\ E(X_{n1}) & \cdots & E(X_{np}) \end{bmatrix}.$$

- ② For a $p \times 1$ random vector \mathbf{x} with mean $\boldsymbol{\mu}_X$ and a $q \times 1$ random vector \mathbf{y} with mean $\boldsymbol{\mu}_Y$, its **covariance matrix** is

$$\text{Cov}(\mathbf{x}, \mathbf{y}) = E[(\mathbf{x} - \boldsymbol{\mu}_X)(\mathbf{y} - \boldsymbol{\mu}_Y)^T] = E(\mathbf{x}\mathbf{y}^T) - \boldsymbol{\mu}_X\boldsymbol{\mu}_Y^T,$$

where its (i, k) th element is $E[(X_i - \mu_{X,i})(Y_k - \mu_{Y,k})]$.

- $\text{Var}(\mathbf{x}) = \text{Cov}(\mathbf{x}, \mathbf{x})$, which is symmetric and positive semi-definite.

Basic Rules: Linear Combination

- ① Random variable X , constants c and d :

$$\begin{aligned}E(cX + d) &= cE(X) + d, \\ \text{Var}(cX + d) &= c^2 \text{Var}(X).\end{aligned}$$

- ② Random vector \mathbf{x} , constant matrix \mathbf{C} , constant vector \mathbf{d} :

$$\begin{aligned}E(\mathbf{C}\mathbf{x} + \mathbf{d}) &= \mathbf{C}E(\mathbf{x}) + \mathbf{d}, \\ \text{Var}(\mathbf{C}\mathbf{x} + \mathbf{d}) &= \mathbf{C}\text{Var}(\mathbf{x})\mathbf{C}^T.\end{aligned}$$

- ③ Random matrix \mathbf{X} , and constant matrices \mathbf{A} , \mathbf{B} , and \mathbf{D}

$$E(\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{D}) = \mathbf{A}E(\mathbf{X})\mathbf{B} + \mathbf{D}.$$

Basic Rules: Additive and Scaling Property

- ① Random variables X , Y , and Z , constant a :

$$\begin{aligned}E(X + Y) &= E(X) + E(Y), \\ \text{Cov}(X + Y, Z) &= \text{Cov}(X, Z) + \text{Cov}(Y, Z), \\ \text{Cov}(aX, Z) &= a\text{Cov}(X, Z).\end{aligned}$$

- ② Random vectors \mathbf{x} , \mathbf{y} , and \mathbf{z} , constant matrices \mathbf{A} and \mathbf{B} :

$$\begin{aligned}E(\mathbf{x} + \mathbf{y}) &= E(\mathbf{x}) + E(\mathbf{y}), \\ \text{Cov}(\mathbf{x} + \mathbf{y}, \mathbf{z}) &= \text{Cov}(\mathbf{x}, \mathbf{z}) + \text{Cov}(\mathbf{y}, \mathbf{z}), \\ \text{Cov}(\mathbf{A}\mathbf{x}, \mathbf{B}\mathbf{y}) &= \mathbf{A}\text{Cov}(\mathbf{x}, \mathbf{y})\mathbf{B}^T.\end{aligned}$$

- ③ Random matrices \mathbf{X} and \mathbf{Y} :

$$E(\mathbf{X} + \mathbf{Y}) = E(\mathbf{X}) + E(\mathbf{Y}).$$

Multivariate Normal Distribution

Let $\mathbf{Z} = [Z_1 \ Z_2 \ \cdots \ Z_p]^T$ be a random vector, each $Z_j \sim N(0, 1)$, and Z_j is independent of Z_k for any $j \neq k$. Let \mathbf{A} be a constant matrix and $\boldsymbol{\mu}$ be a constant vector. Then,

$$\mathbf{X} = \mathbf{AZ} + \boldsymbol{\mu}$$

follows a p -dimensional **multivariate normal distribution**. We often denote a multivariate normal random vector by $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu}$ is the mean and $\boldsymbol{\Sigma}$ is the covariance matrix.

Using multivariate normal distribution, we can write the assumption $e | X = x \sim N(0, \sigma^2)$ as

$$e | \mathbf{x} \sim N(0, \sigma^2 \mathbf{I}).$$

Properties Needed In Our Course

- ① If $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, \mathbf{A} is a matrix of constants, and \mathbf{d} is a vector of constants, then $\mathbf{AX} + \mathbf{d} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{d}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$.
- ② Let $\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$ be distributed as $N\left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}\right)$.
 - ① All subsets of \mathbf{x} are normally distributed, i.e., \mathbf{X}_1 is distributed as $N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$, and \mathbf{X}_2 is distributed as $N(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$.
 - ② If X_i is the i th element in \mathbf{x} , then $X_i \sim N(\mu_i, \Sigma_{ii})$.
 - ③ The conditional distribution of \mathbf{x}_1 given that \mathbf{x}_2 , is

$$\mathbf{x}_1 | \mathbf{x}_2 \sim N\left\{\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}\right\},$$

provided that $\boldsymbol{\Sigma}_{22}$ is invertible.

Distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$

Under the normality assumption, we can obtain

$$\begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} \sim N \left(\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}, \begin{bmatrix} \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] & -\frac{\sigma^2 \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ -\frac{\sigma^2 \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} & \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{bmatrix} \right).$$

The marginal distributions are

$$\begin{aligned} \hat{\beta}_0 &\sim N \left(\beta_0, \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \right), \\ \hat{\beta}_1 &\sim N \left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right). \end{aligned}$$

The **standard errors** of $\hat{\beta}_0$ and $\hat{\beta}_1$ are

$$\hat{\sigma} \sqrt{\left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]} \quad \text{and} \quad \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}, \text{ respectively.}$$

Student t-Distribution

It can be shown that, conditional on \mathbf{X} ,

$$\frac{\sum_{i=1}^n \hat{e}_i^2}{\sigma^2} \sim \chi^2(n-2).$$

Then,

$$\frac{(\hat{\beta}_0 - \beta_0) / \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}}{\sqrt{\sum_{i=1}^n \hat{e}_i^2 / (n-2)}} \sim t(n-2),$$

and

$$\frac{(\hat{\beta}_1 - \beta_1) / \sqrt{\frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2}}}{\sqrt{\sum_{i=1}^n \hat{e}_i^2 / (n-2)}} \sim t(n-2).$$

Confidence Interval

A $1 - \alpha$ confidence interval for β_0 is

$$\hat{\beta}_0 \pm t_{1-\alpha/2}(n-2) \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)},$$

where $t_{1-\alpha/2}(n-2)$ satisfies

$$P(T \leq t_{1-\alpha/2}(n-2)) = 1 - \frac{\alpha}{2},$$

with $T \sim t(n-2)$.

A $1 - \alpha$ confidence interval for β_1 is

$$\hat{\beta}_1 \pm t_{1-\alpha/2}(n-2) \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$