

# Linear Regression

# Learning Objectives

- Matrix Formulation
- Multiple Linear Regression
- Model Assumptions

# Matrix Formulation

# Matrix Version of Model

$$Y_i = \beta_0 + \beta_1 X_{1i} + \epsilon_i$$

$$Y_i = \mathbf{X}_i^T \boldsymbol{\beta} + \epsilon_i$$

- $Y_i$ : Outcome Variable
- $\mathbf{X}_i = (1, X_i)^T$ : Predictors
- $\boldsymbol{\beta} = (\beta_0, \beta_1)^T$ : Coefficients
- $\epsilon_i$ : error term

$$Y_i = \begin{matrix} \mathbf{X}_i^T \\ 1 \times 2 \end{matrix} \begin{matrix} \boldsymbol{\beta} \\ 2 \times 1 \end{matrix} + \begin{matrix} \epsilon_i \\ 1 \times 1 \end{matrix}$$

# Data Matrix Formulation


$$X_i = (1 \quad X_i)^T$$

For  $n$  data points

$$\begin{matrix} \mathbf{Y} & = & \mathbf{X}^T \boldsymbol{\beta} & + & \boldsymbol{\epsilon} \\ \uparrow_{n \times 1} & & \uparrow_{n \times 2} \uparrow_{2 \times 1} & & \uparrow_{n \times 1} & & = & \uparrow_{n \times 1} \end{matrix}$$

- $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ : Outcome Variable
- $\mathbf{X} = (\overset{p \times 1}{X_1}, \dots, X_n)^T$ : Predictors
- $\overset{p \times n}{\boldsymbol{\beta}} = (\beta_0, \beta_1)^T$ : Coefficients
- $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$ : Error terms

# Least Squares Formula

$$(Y - X^T \beta)^T (Y - X^T \beta)$$

$$\sum_{i=1}^n (Y_i - (\beta_0 + \beta_1 X_i))^2$$

# Estimates

$$\hat{\beta} = \overset{2 \times n}{(X^T X)^{-1}} \overset{n \times 1}{X^T Y} \quad 2 \times 1$$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \begin{pmatrix} \bar{y} - \hat{\beta}_1 \bar{x} \\ \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \end{pmatrix}$$

# Multiple Linear Regression



# MLR

Multivariable linear regression models are used when more than one explanatory variable is used to explain the outcome of interest.

# Continuous Variable

$p+1 \ll N$

To fit an additional continuous random variable to the model, we will only need to add it to the model:

100 99

overfitting.

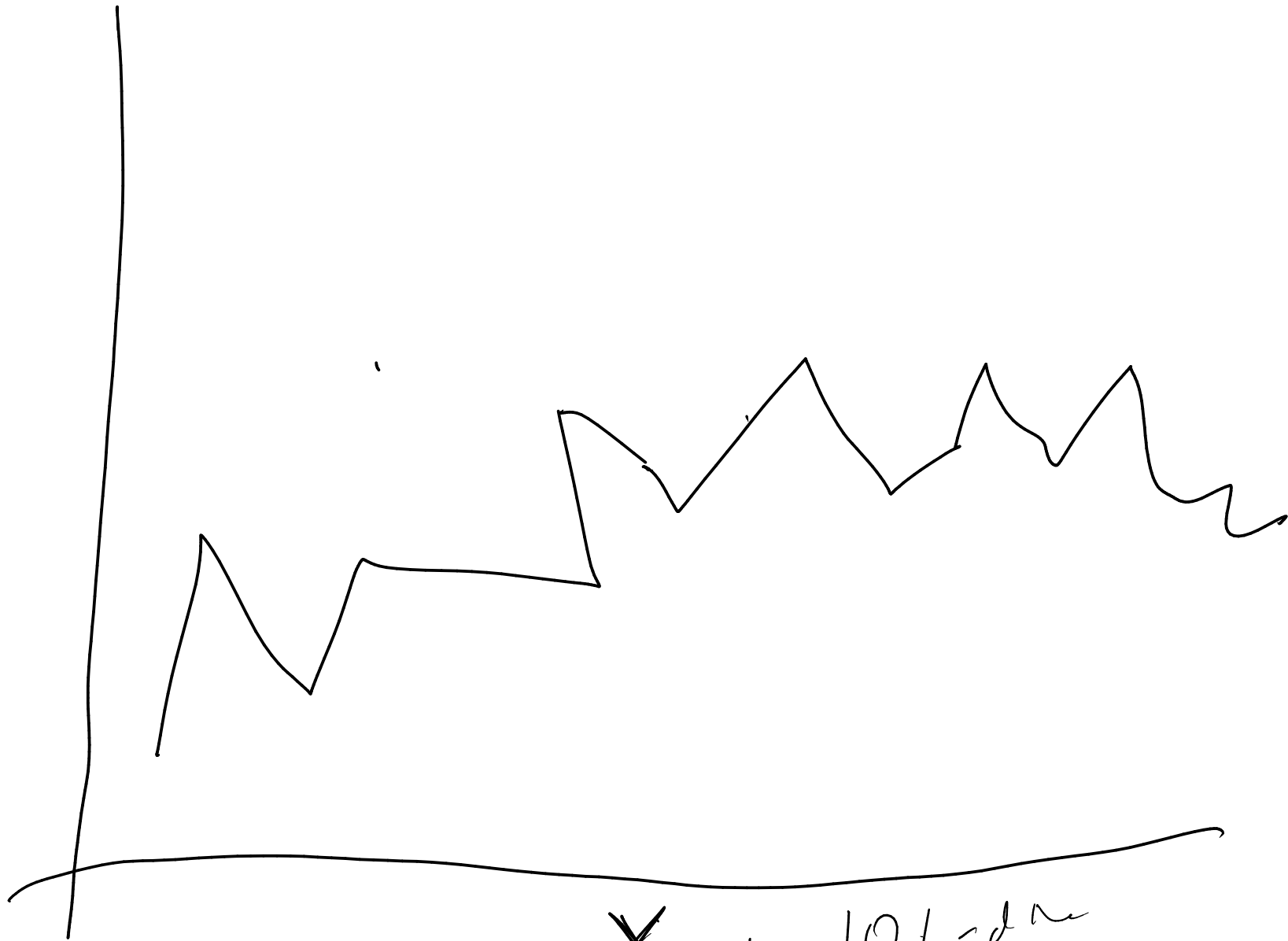
$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

body mass      wingspan      Beak size

As  $X_1$  increase by 1 unit,  $Y$  will

increase/decrease by an average of  $\beta_1$

units, adjusting for Beak size



$x_{n+1}$   $10^{-d} n$

# Categorical Variable

$$\text{Species} = \begin{cases} 1 & \text{seagull} \\ 2 & \text{condor} \\ 3 & \text{duck} \end{cases}$$

A categorical variable can be included in a model, but a reference category must be specified.

$$Y = \beta_0 + \beta_1 (\text{wins}) + \beta_2 D_2 + \beta_3 D_3$$

$$\text{Species} \quad D_2 = \begin{cases} 1 & \text{condor} \\ 0 & \text{o.w.} \end{cases} \quad D_3 = \begin{cases} 1 & \text{duck} \\ 0 & \text{o.w.} \end{cases}$$

reference category      seagull

# Fitting a model with categorical variables

To fit a model with categorical variables, we must utilize dummy (binary) variables that indicate which category is being referenced. We use  $C - 1$  dummy variables where  $C$  indicates the number of categories. When coded correctly, each category will be represented by a combination of dummy variables.

# Example

If we have 4 categories, we will need 3 dummy variables:

	Cat 1	Cat 2	Cat 3	Cat 4
Dummy 1	1	0	0	0
Dummy 2	0	1	0	0
Dummy 3	0	0	1	0

Which one is the reference category?

Cat. 4

# Matrix Notation

$$Y_i = \beta^T X_i$$

- $\beta$ : a column vector of regression coefficients
- $X$ : a column vector of predictor variables

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$$

~~$X$~~  =  $\begin{pmatrix} | \\ \text{Wingspan} \\ \text{Are they condor?} \\ \text{Are they duck?} \end{pmatrix}$

$$\hat{\beta} = \begin{pmatrix} X & X^T \end{pmatrix}^{-1} \begin{pmatrix} X & Y \end{pmatrix} = \begin{pmatrix} \beta_0 \\ \hat{\beta}_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix}$$

$4 \times n$     $n \times 4$     $4 \times n$     $n \times 1$     $4 \times 1$



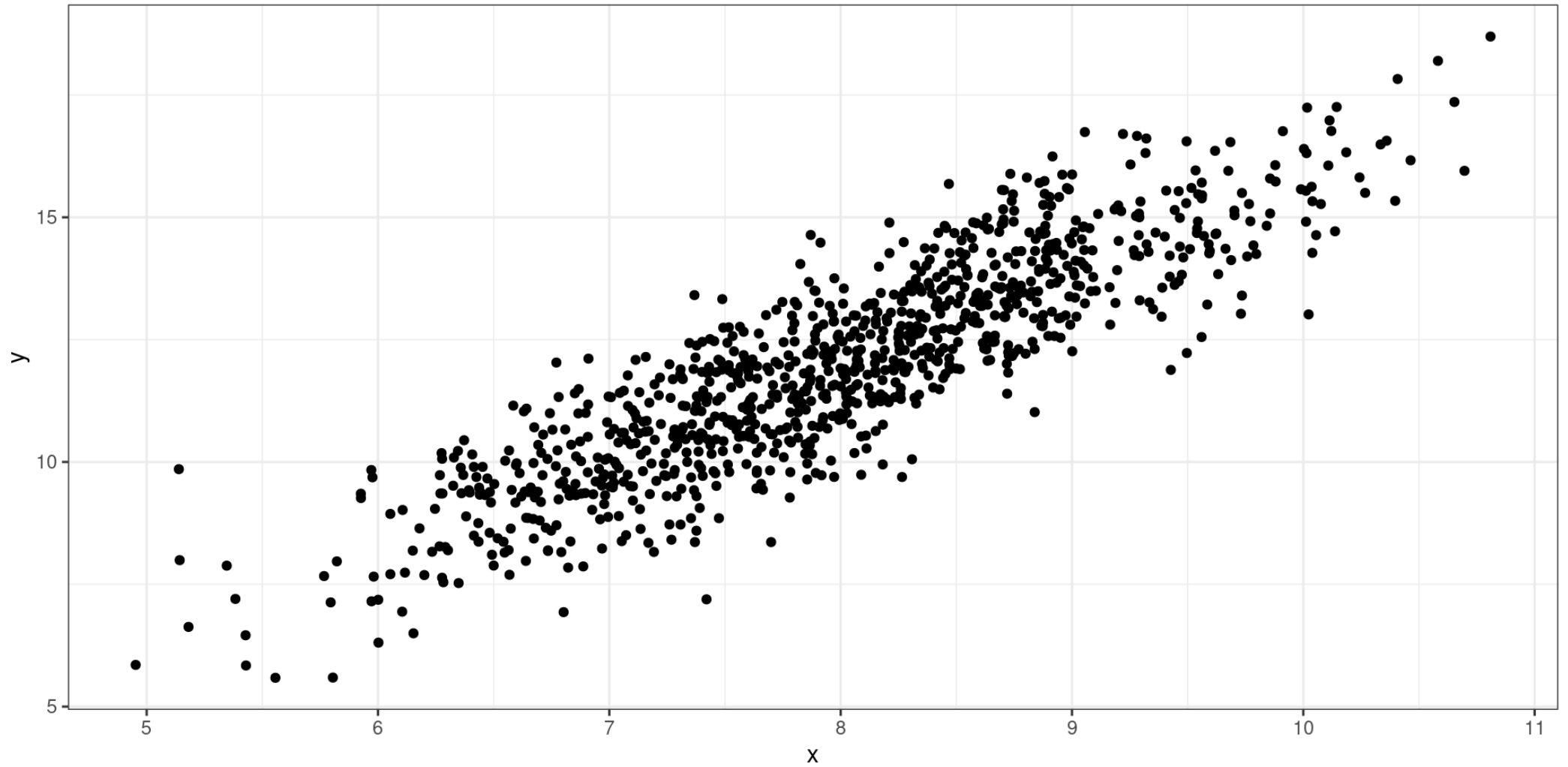
# Model Assumptions

# Model

$$Y = \boldsymbol{\beta}^T \mathbf{X} + \epsilon$$

- $\epsilon \sim N(0, \sigma^2)$

# Model Scatter Plot



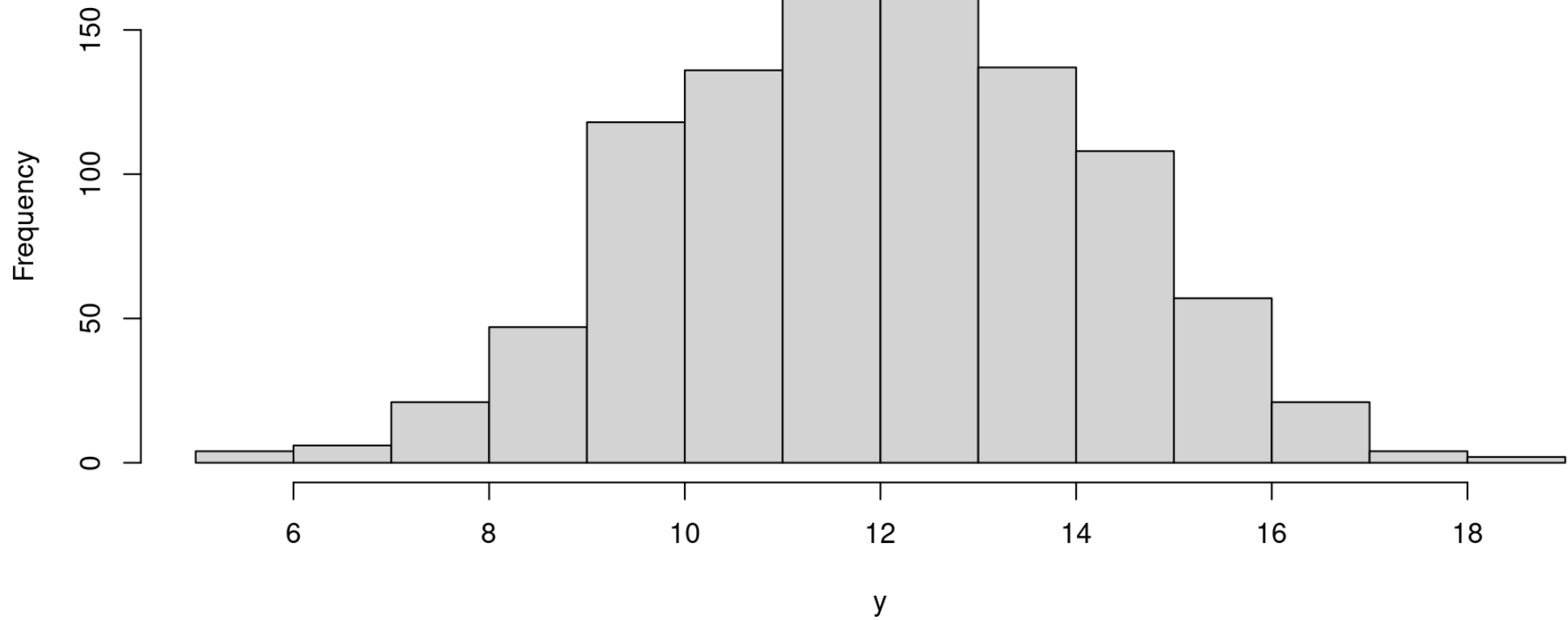
# Model Assumptions

- Errors are normally distributed
- Constant Variance
- Linearity
- Independence
- No outliers

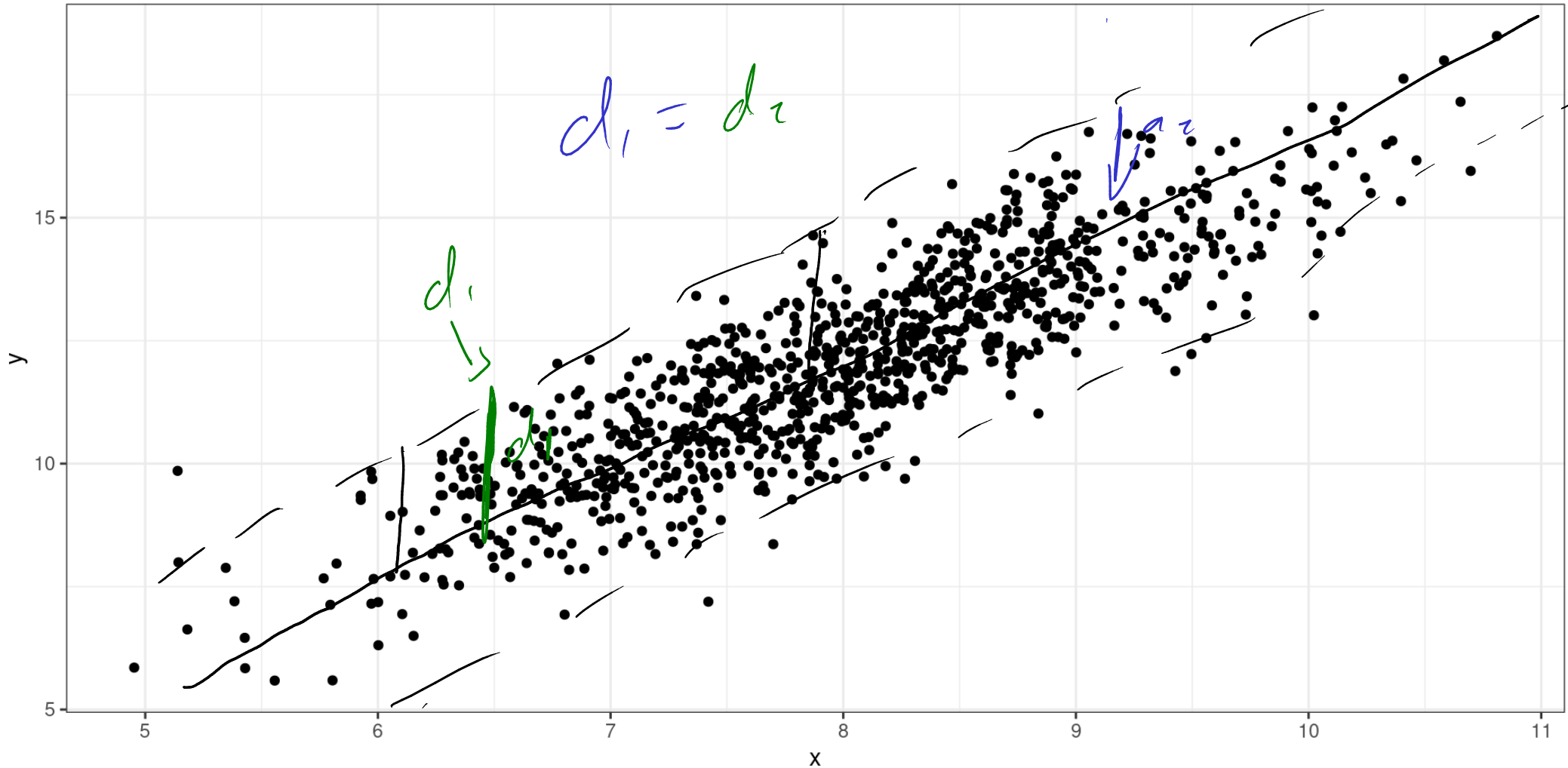
$$v_i = y_i - \hat{y}$$

# Errors Normally Distributed

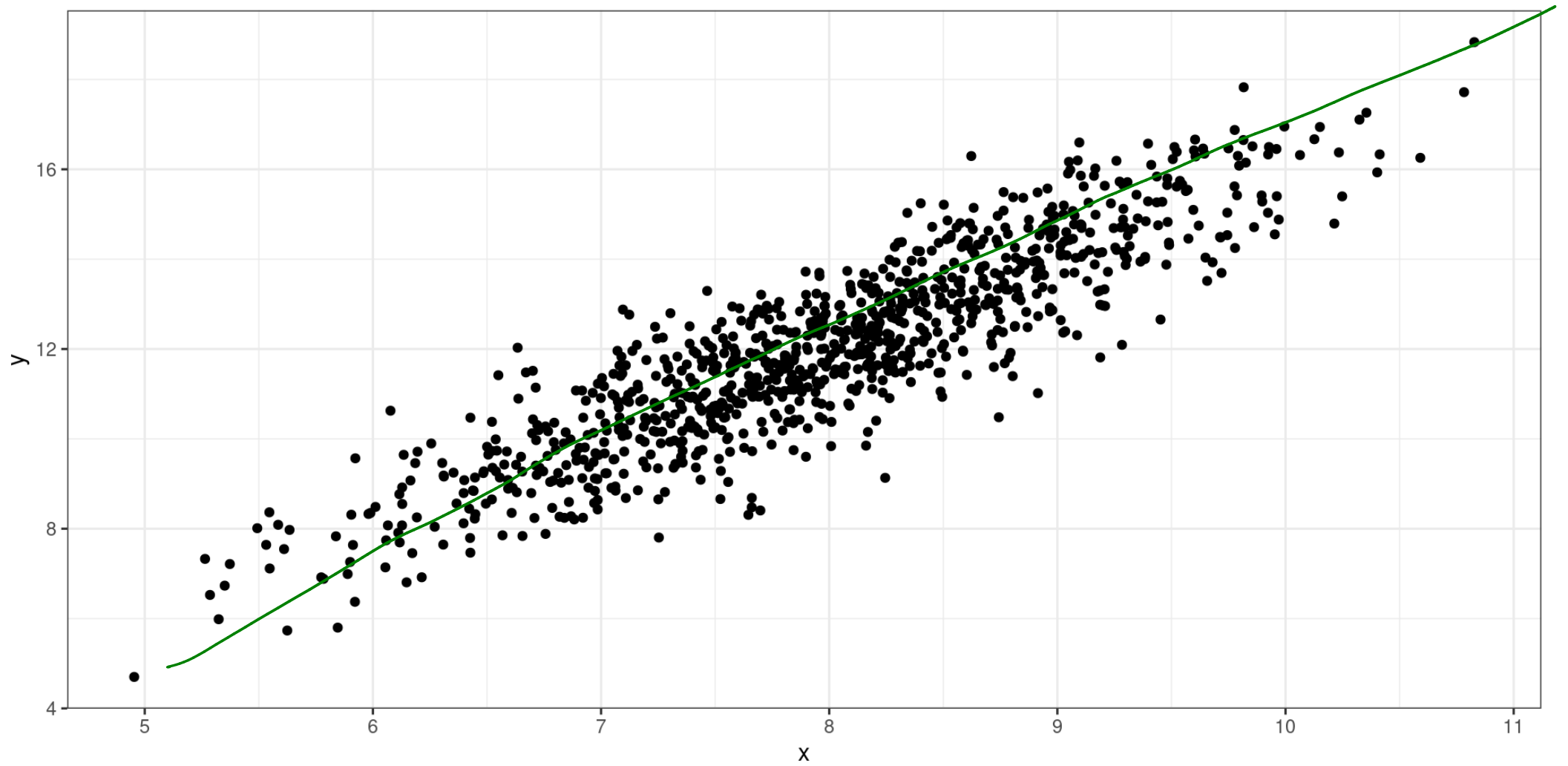
Histogram of Error



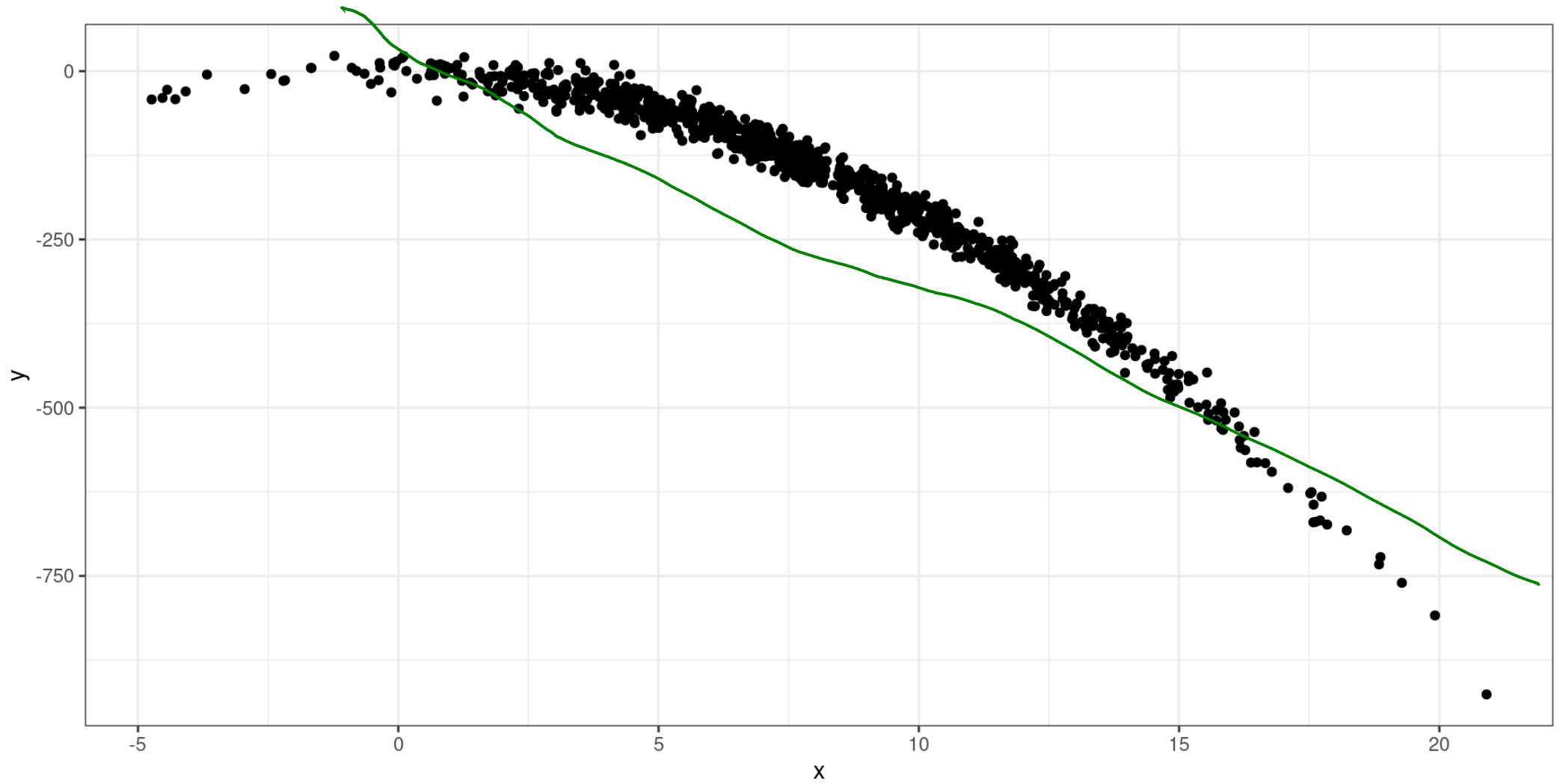
# Constant Variance



# Linearity

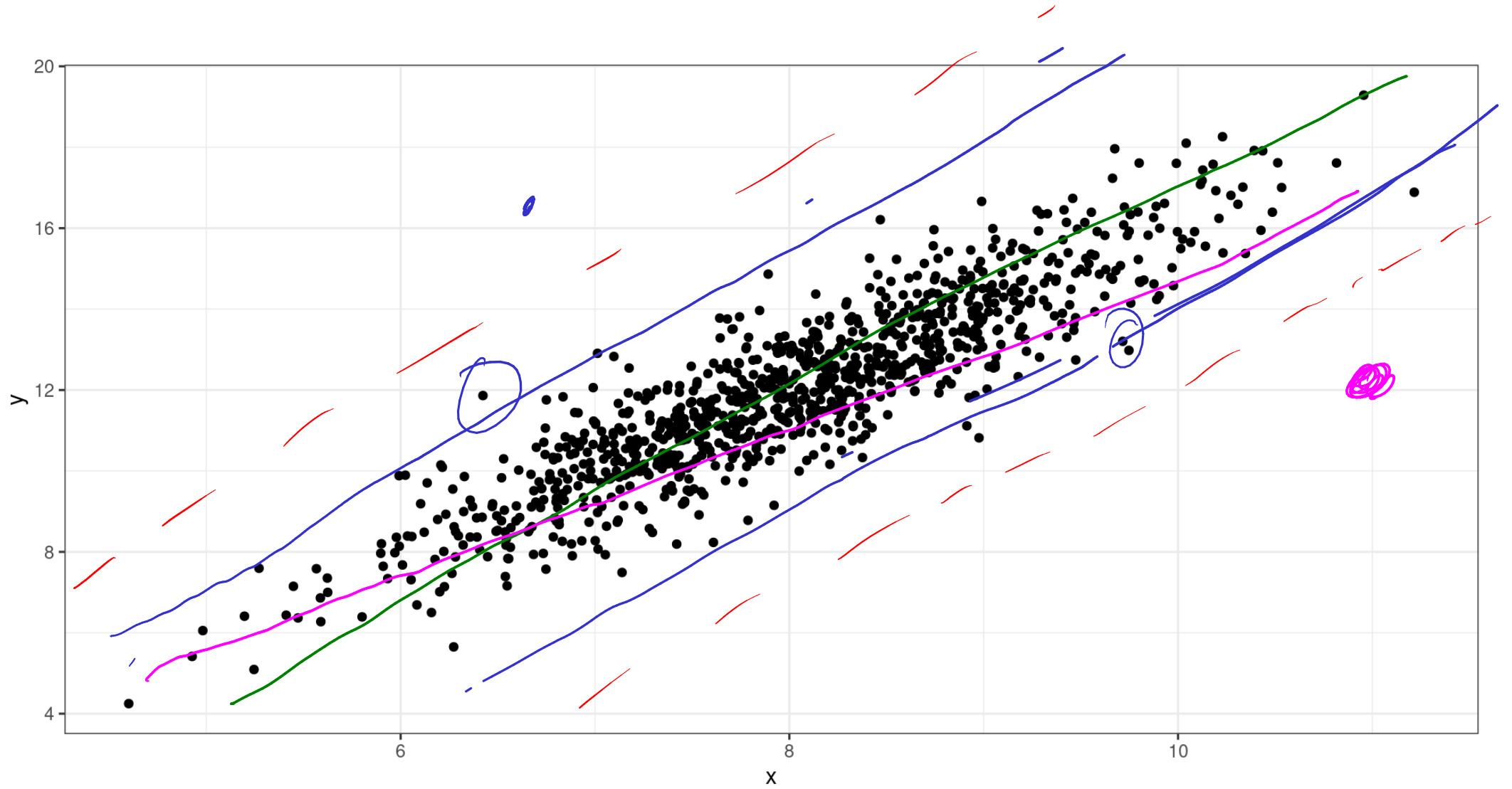


# Linearity





# No Outliers



# Residual Analysis

A residual analysis is used to assess the validity of the assumptions.

$$\text{Residuals} = Y_i - \hat{Y}_i$$

Standardized

Studentized

Hat values

Cook's Distance

