

Linear Regression

Learning Objectives

- Matrix Formulation
- Multiple Linear Regression
- Model Assumptions

Matrix Formulation

Matrix Version of Model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$\underset{1 \times 1}{Y_i} = \underset{1 \times 2}{X_i^T} \underset{2 \times 1}{\beta} + \underset{1 \times 1}{\epsilon_i}$$

$$(x_1 \quad x_2) \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$

$$x_1 \beta_0 + x_2 \beta_1$$

- Y_i : Outcome Variable
- $X_i = (1, X_i)^T$: Predictors
- $\beta = (\beta_0, \beta_1)^T$: Coefficients
- ϵ_i : error term

Data Matrix Formulation

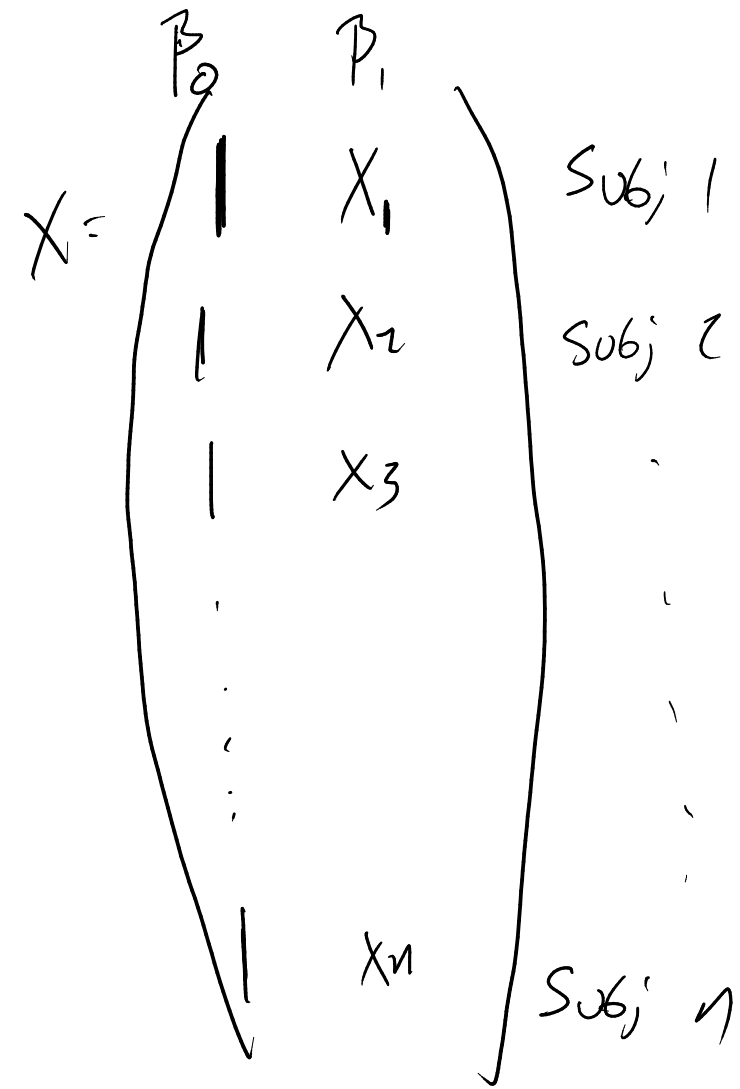
For n data points

$$Y_i = \beta_0 + \beta_1 X_i$$


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

$n \times 1$ $n \times 2$ 2×1 $n \times 1$

- $\mathbf{Y} = (Y_1, \dots, Y_n)^T$: Outcome Variable
- $\mathbf{X} = (X_1, \dots, X_n)^T$: Predictors
- $\boldsymbol{\beta} = (\beta_0, \beta_1)^T$: 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

$$\underset{n \times 2}{\mathbf{X}} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$

$$\mathbf{Y} = \begin{pmatrix} 1 & 1 & 1 & \dots & 1 \\ x_1 & \dots & \dots & \dots & x_n \end{pmatrix}$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

$2 \times n \quad n \times 2 \quad 2 \times n \quad n \times 1$

$$\hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix}$$

2×1

Multiple Linear Regression

MLR

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

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

wing *beak*

As wingspan increases by 1 unit, body mass will increase/decrease by an average of β_1 units, after adjusting for beak size

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}))^2$$

$$X = \begin{pmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} \end{pmatrix}$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

3×1 $3 \times n$ $n \times 3$ $3 \times n$ $n \times 1$

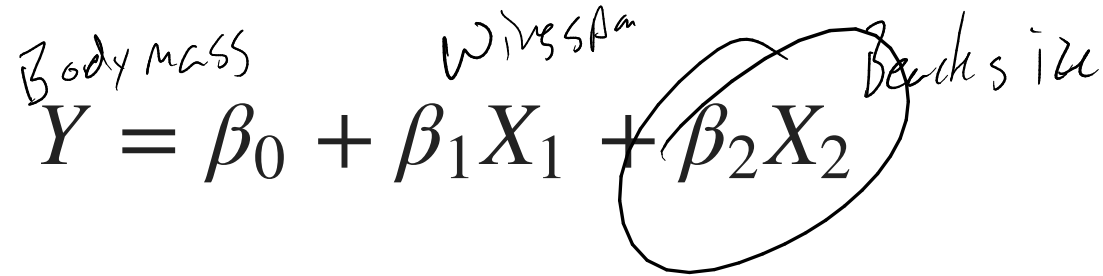
$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}$$

Continuous Variable

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

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

Body mass Wires size Beach size



Categorical Variable

$$X_2 = \begin{cases} 1 & \text{Seagull} \\ 2 & \text{condor} \\ 3 & \text{duck} \end{cases}$$

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

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 D_2 + \beta_3 D_3$$

$\hat{Y}_{\text{duck}} - \hat{Y}_{\text{condor}} = \hat{\beta}_0 + \hat{\beta}_1(2) + \beta_2 - (\hat{\beta}_0 + \hat{\beta}_1(1) + \beta_2)$
 $\hat{\beta}_1(2-1) + \beta_3 - \beta_2$

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

Reference category = Seagull

On average the body mass of a ~~condor~~ duck is greater/

lesser than a seagull by Bg units, adjusting for wingspan

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

$$Y = \beta^T X$$

- β : a column vector of regression coefficients
- X : a column vector of predictor variables

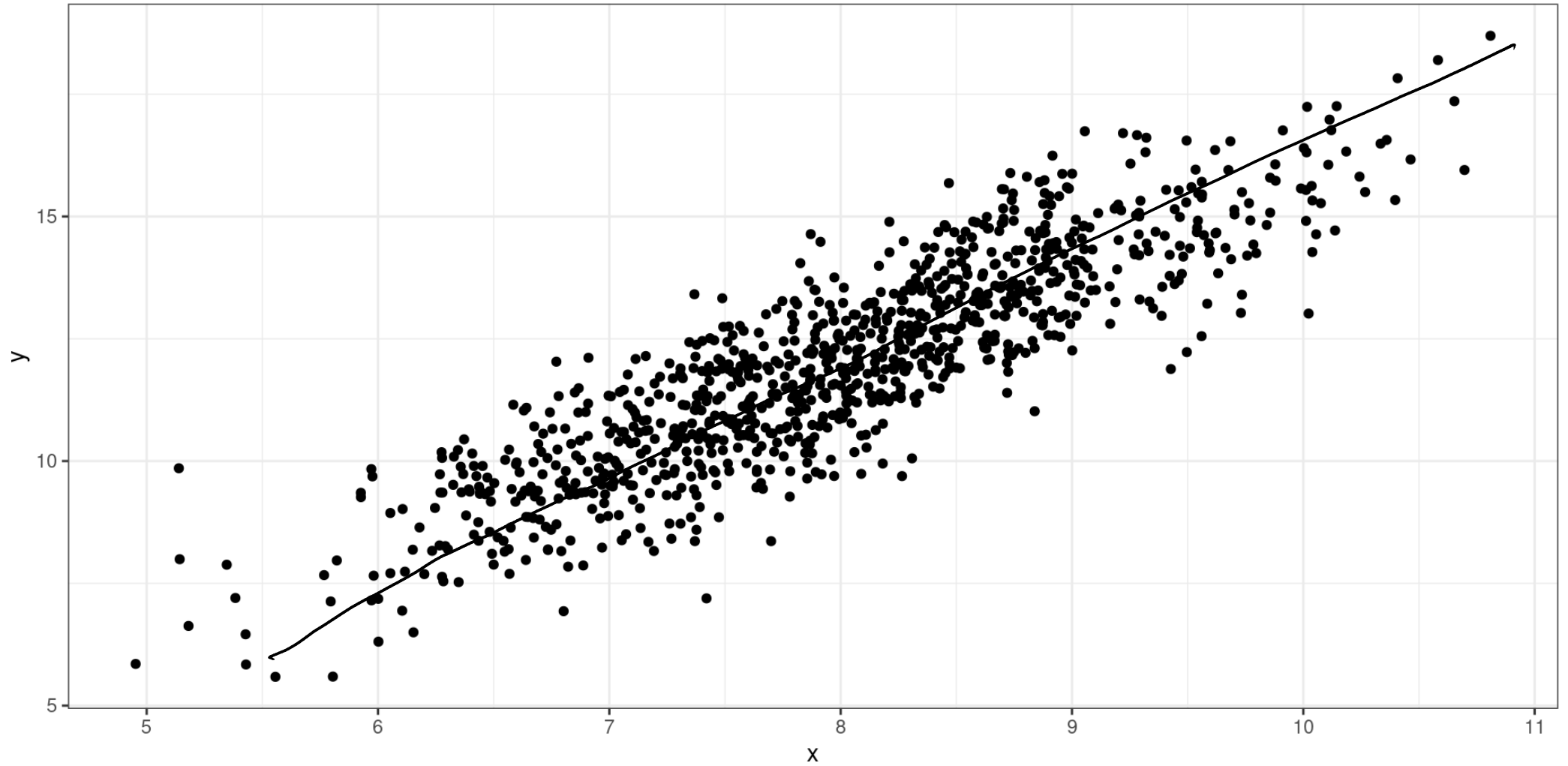
Model Assumptions

Model

$$Y = \beta^T X$$

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

Model Scatter Plot



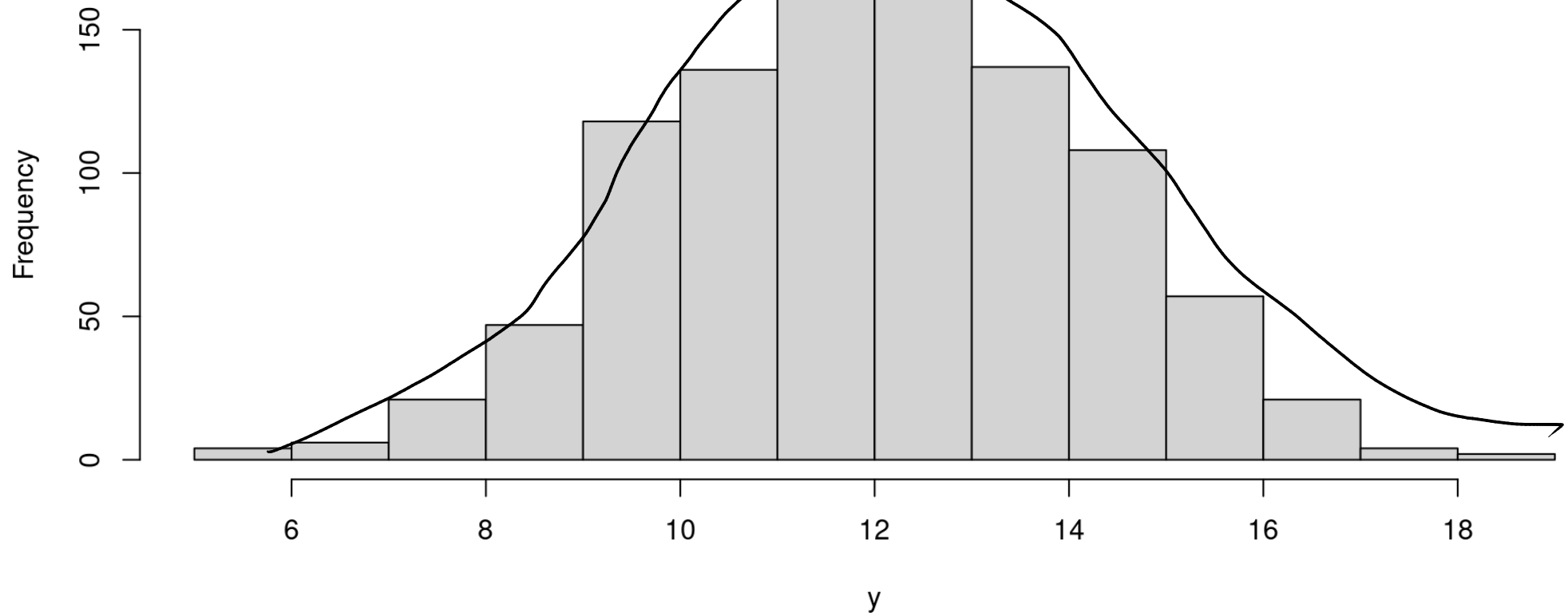
(x_1, x_2, x_3)

Model Assumptions

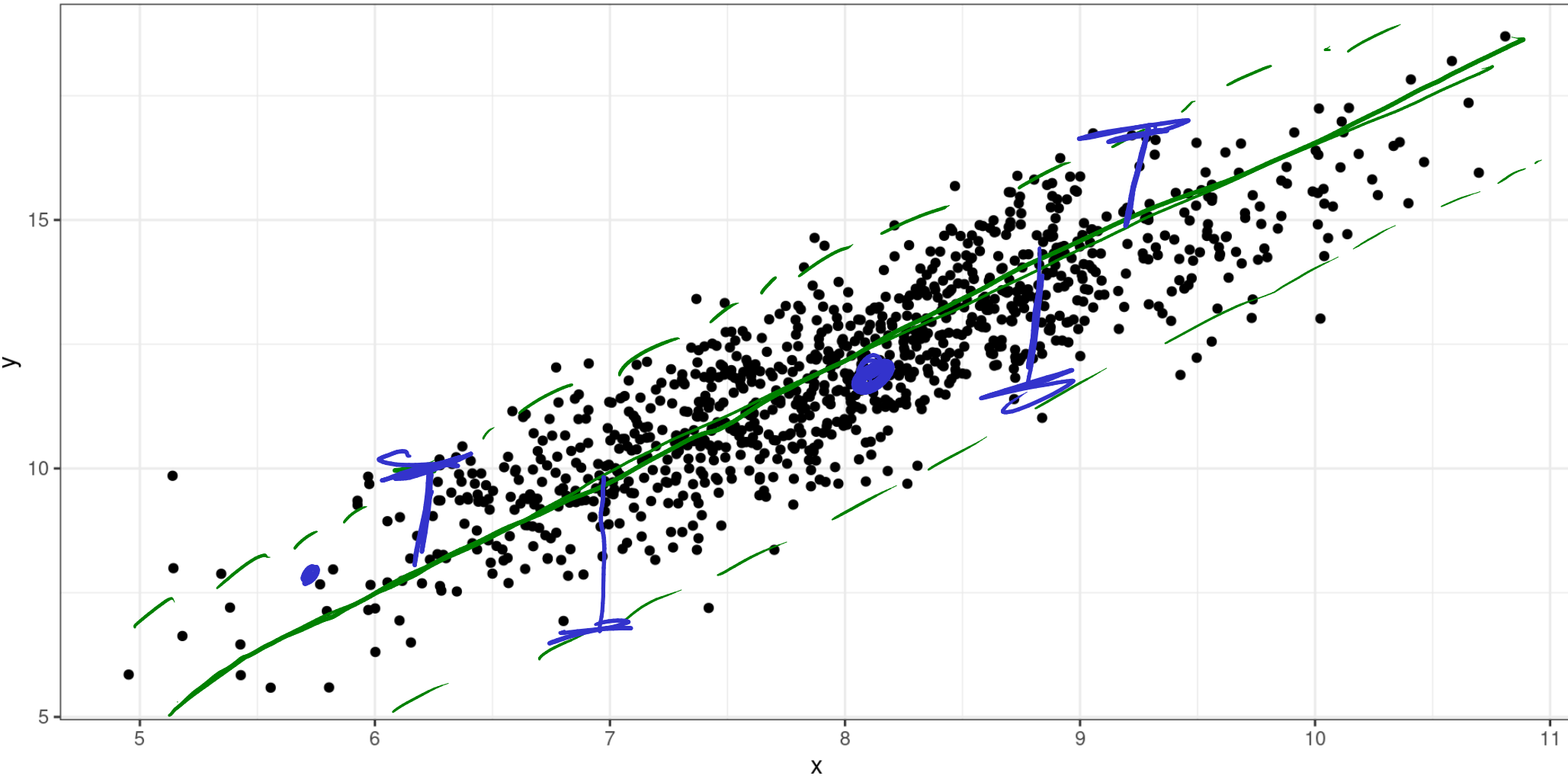
- Errors are normally distributed
- Constant Variance
- Linearity
- Independence
- No outliers

Errors Normally Distributed

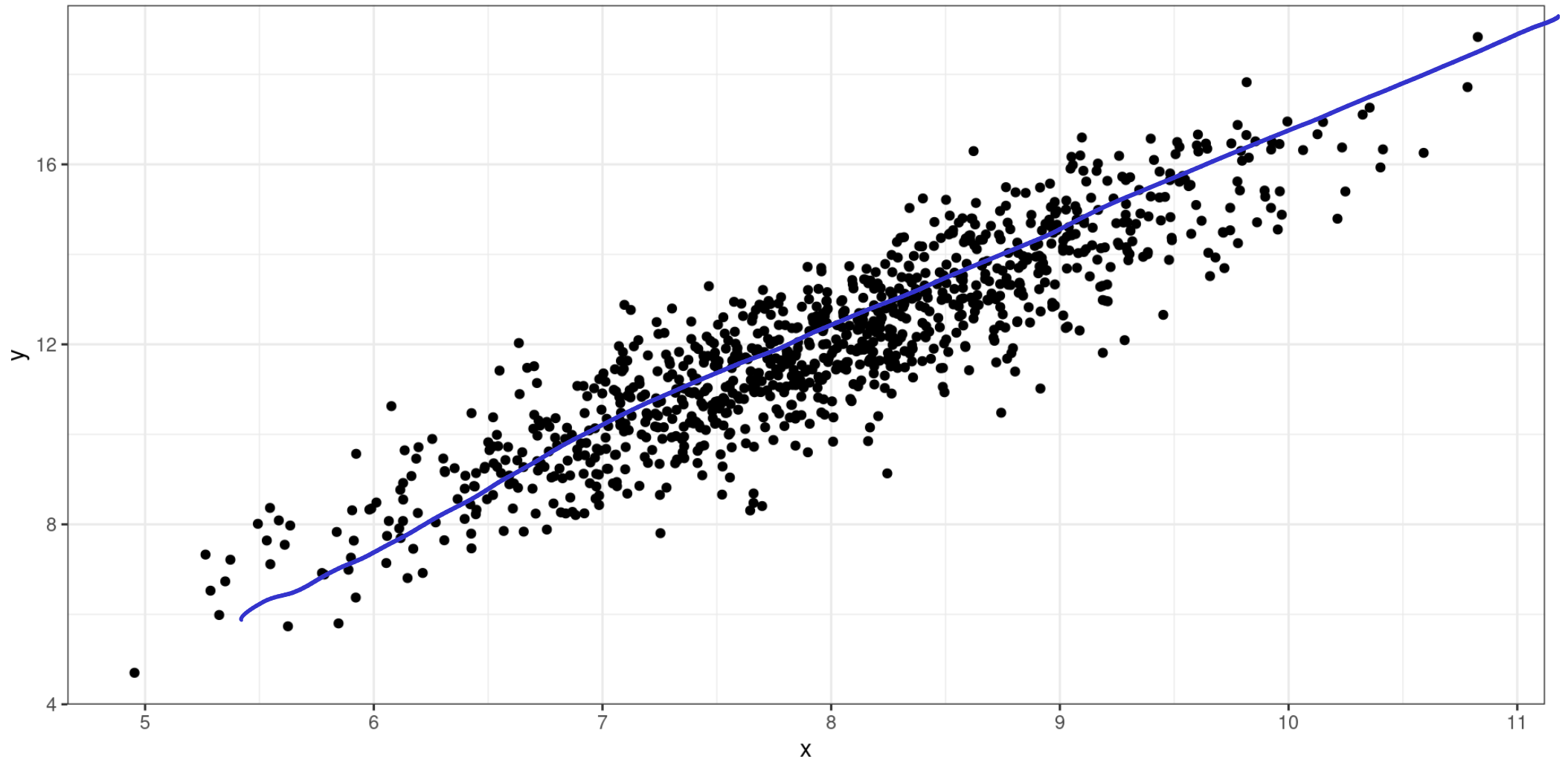
Histogram of Error



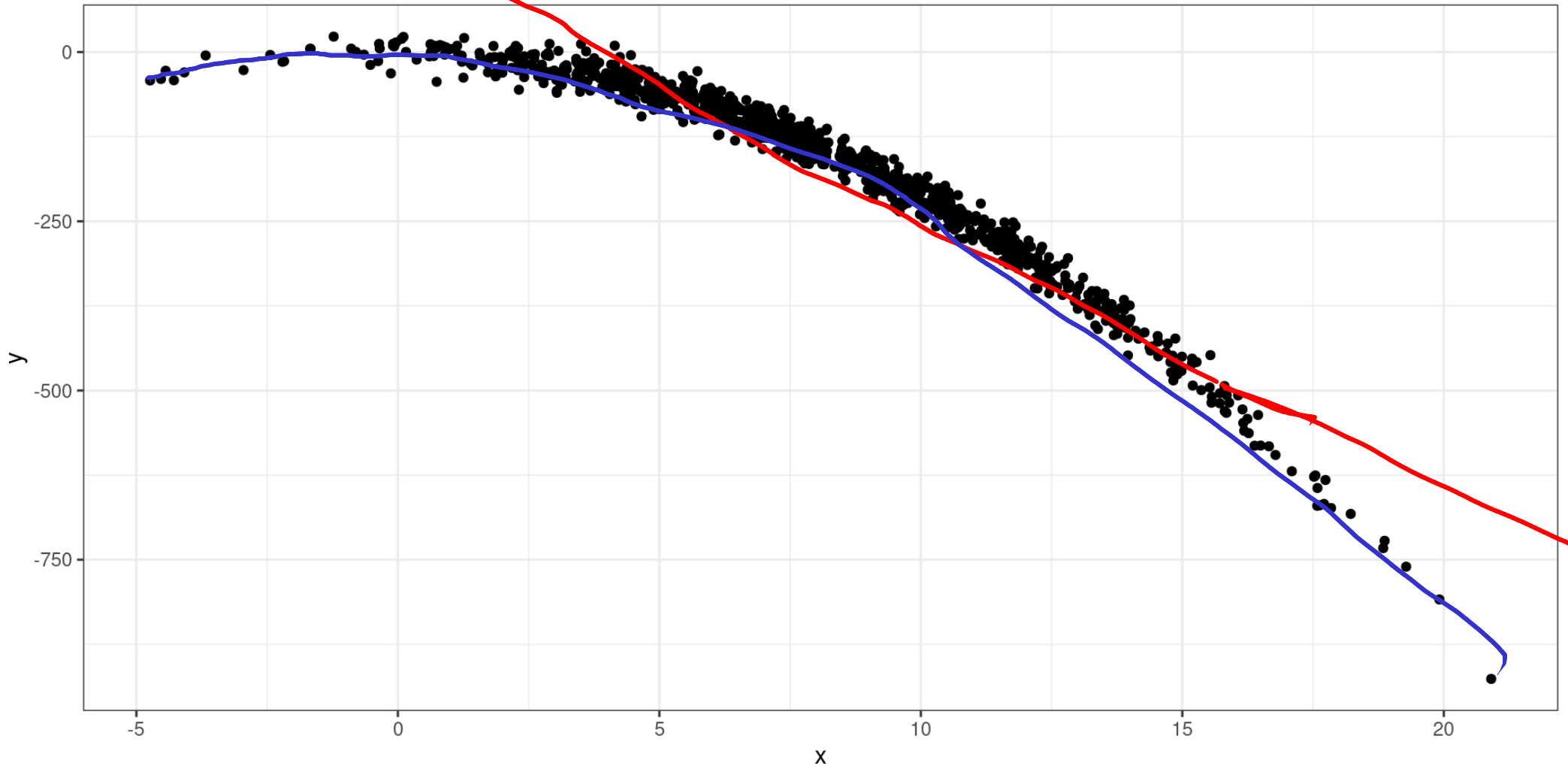
Constant Variance



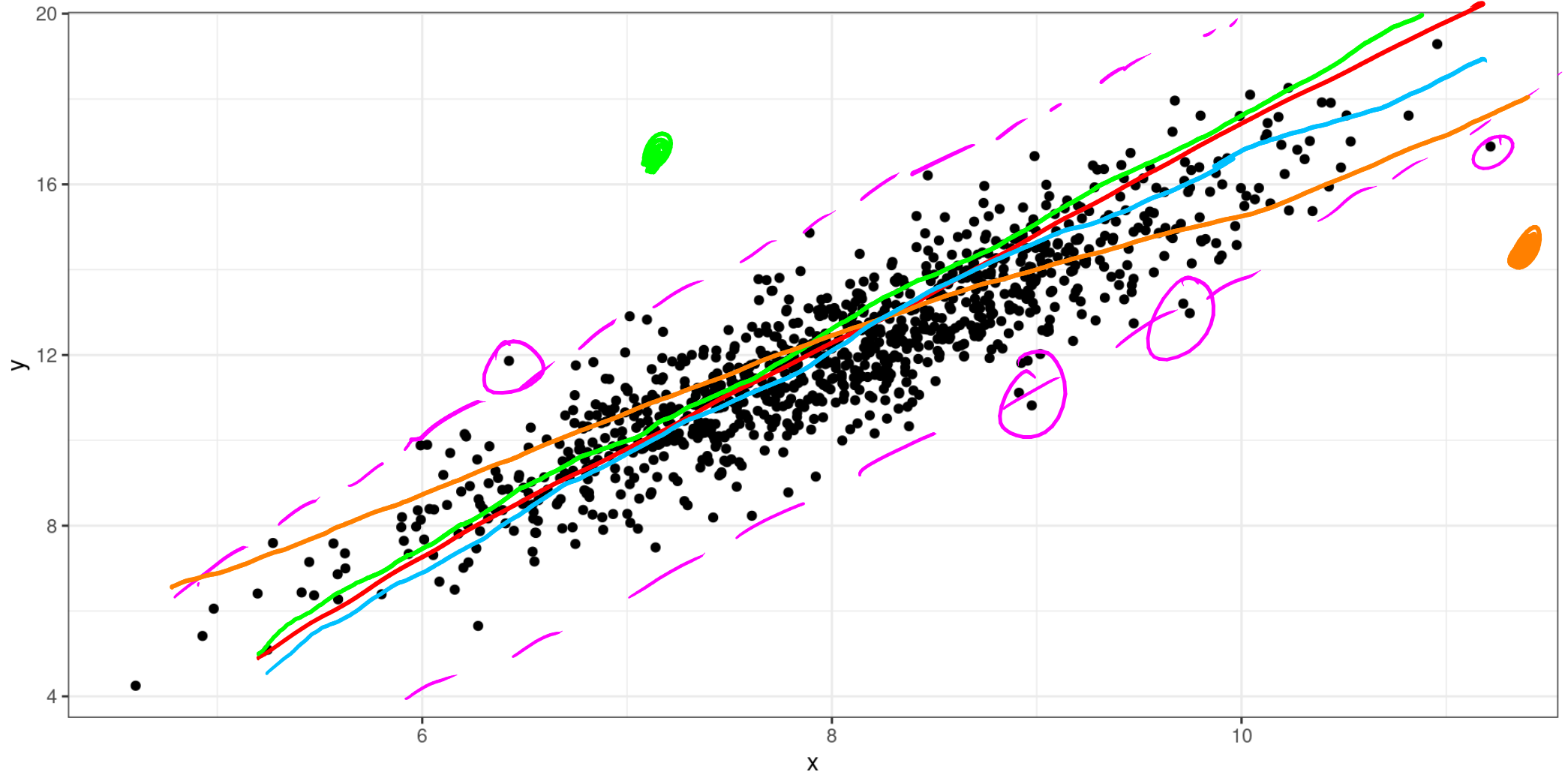
Linearity



Linearity



No Outliers



Residual Analysis

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

$$r_i = y_i - \hat{y}_i$$



Studentized
Standardized
hat values
Cook's Distance

