

Regression Analysis

The Basics

$$\mathbf{Y} = \beta_0 + \beta_1 x + \epsilon, \quad \hat{\beta}_1 = \frac{S_{XY}}{S_{XX}}, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = \bar{Y} - \frac{S_{XY}}{S_{XX}} \bar{X},$$

$$\text{where } S_{XX} = \sum_i (\mathbf{X}_i - \bar{\mathbf{X}})^2 \quad \text{and} \quad S_{XY} = \sum_i (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{Y}_i - \bar{\mathbf{Y}}).$$

Fixed Design: If x_i is fixed, $\sum_i (x_i - \bar{x}) = 0$ and therefore, $S_{XX} = \sum_i (\mathbf{X}_i - \bar{\mathbf{X}})^2 = \sum_i (x_i - \bar{x})^2 = S_{xx}$,

$$\text{and } S_{XY} = \sum_i (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{Y}_i - \bar{\mathbf{Y}}) = \sum_i (x_i - \bar{x})(\mathbf{Y}_i - \bar{\mathbf{Y}}) = \sum_i (x_i - \bar{x})\mathbf{Y}_i - \bar{\mathbf{Y}} \sum_i (x_i - \bar{x}) = \sum_i (x_i - \bar{x})\mathbf{Y}_i.$$

The Gradient

$$\begin{aligned} \mathbf{E}(\hat{\beta}_1) &= \mathbf{E}\left(\frac{S_{XY}}{S_{XX}}\right) = \frac{\mathbf{E}(S_{XY})}{S_{xx}} = \frac{\mathbf{E}(\sum_i (x_i - \bar{x})\mathbf{Y}_i)}{S_{xx}} = \frac{\sum_i (x_i - \bar{x})\mathbf{E}(\mathbf{Y}_i)}{S_{xx}} = \frac{\sum_i (x_i - \bar{x})\mathbf{E}(\beta_0 + \beta_1 x_i)}{S_{xx}} \\ &= \frac{\beta_0}{S_{xx}} \sum_i (x_i - \bar{x}) + \frac{\beta_1}{S_{xx}} \sum_i x_i (x_i - \bar{x}) = \frac{\beta_1}{S_{xx}} \left(\sum_i x_i (x_i - \bar{x}) - \sum_i \bar{x} (x_i - \bar{x}) \right) = \frac{\beta_1}{S_{xx}} \sum_i (x_i - \bar{x})^2 = \beta_1. \end{aligned}$$

$$\begin{aligned} \mathbf{Var}(\hat{\beta}_1) &= \mathbf{Var}\left(\frac{S_{XY}}{S_{XX}}\right) = \mathbf{Var}\left(\frac{\sum_i (x_i - \bar{x})\mathbf{Y}_i}{S_{xx}}\right) = \mathbf{Var}\left(\sum_i \frac{x_i - \bar{x}}{S_{xx}} \mathbf{Y}_i\right) = \sum_i \left(\frac{x_i - \bar{x}}{S_{xx}}\right)^2 \mathbf{Var}(\mathbf{Y}_i) \\ &= \frac{\mathbf{Var}(\mathbf{Y}_i)}{S_{xx}^2} \sum_i (x_i - \bar{x})^2 = \frac{\sigma^2}{S_{xx}^2} S_{xx} = \frac{\sigma^2}{S_{xx}}. \end{aligned}$$

$$\boxed{\hat{\beta}_1 \sim \mathbf{N}\left(\beta_1, \frac{\sigma}{\sqrt{S_{xx}}}\right)} \quad \sqrt{S_{xx}} \cdot \frac{\hat{\beta}_1 - \beta_1}{\sigma} \sim \mathbf{N}(0, 1), \quad \sqrt{S_{xx}} \cdot \frac{\hat{\beta}_1 - \beta_1}{S} \sim t_{n-2}.$$

The Intercept

$$\begin{aligned} \mathbf{E}(\hat{\beta}_0) &= \mathbf{E}(\bar{Y} - \hat{\beta}_1 \bar{X}) = \mathbf{E}(\bar{Y}) - \mathbf{E}(\hat{\beta}_1 \bar{X}) = \mathbf{E}\left(\frac{\sum_i \mathbf{Y}_i}{n}\right) - \mathbf{E}(\hat{\beta}_1 \bar{x}) = \mathbf{E}\left(\frac{\sum_i \beta_0 + \beta_1 x_i}{n}\right) - \bar{x} \mathbf{E}(\hat{\beta}_1) \\ &= \frac{n\beta_0}{n} + \beta_1 \sum_i \frac{x_i}{n} - \bar{x} \beta_1 = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} = \beta_0. \end{aligned}$$

$$\begin{aligned} \mathbf{Var}(\hat{\beta}_0) &= \mathbf{Var}(\bar{Y} - \hat{\beta}_1 \bar{X}) = \mathbf{Var}\left(\frac{\sum_i \mathbf{Y}_i}{n}\right) - \bar{x}^2 \mathbf{Var}(\hat{\beta}_1) = \frac{\sum_i \mathbf{Var}(\mathbf{Y}_i)}{n^2} - \bar{x}^2 \frac{\sigma^2}{S_{xx}} \\ &= \frac{n\sigma^2}{n^2} - \bar{x}^2 \frac{\sigma^2}{S_{xx}} = \sigma^2 \left(\frac{1}{n} - \frac{\bar{x}^2}{S_{xx}}\right). \end{aligned}$$

$$\boxed{\hat{\beta}_0 \sim \mathbf{N}\left(\beta_0, \sigma \sqrt{\frac{1}{n} - \frac{\bar{x}^2}{S_{xx}}}\right)} \quad \frac{\hat{\beta}_0 - \beta_0}{\sigma \sqrt{\frac{1}{n} - \frac{\bar{x}^2}{S_{xx}}}} \sim \mathbf{N}(0, 1), \quad \frac{\hat{\beta}_0 - \beta_0}{S \sqrt{\frac{1}{n} - \frac{\bar{x}^2}{S_{xx}}}} \sim t_{n-2}.$$

Test Statistic

$$t_0 = \sqrt{S_{xx}} \cdot \frac{\hat{b}_1}{S}, \quad \text{where } \hat{b}_1 \text{ is an estimate of } \hat{\beta}_1.$$

To test the null hypotheses $H_0 : \beta_2 = 0$ against the alternative hypothesis $H_a : \beta_2 \leq 0$,

$$H_0 \text{ is to be rejected if } \hat{b}_1 \notin \left[-t \frac{S}{\sqrt{S_{xx}}}, t \frac{S}{\sqrt{S_{xx}}}\right], \quad \text{where } t = t_{n-1, 1-\frac{\alpha}{2}}. \quad \text{i.e. } |\hat{b}_1| > t \frac{S}{\sqrt{S_{xx}}} = \left|\frac{\hat{b}_1}{t_0}\right|.$$

$$\boxed{\text{Reject } H_0 \text{ if } |t_0| > t_{n-1, 1-\frac{\alpha}{2}}.} \quad \text{i.e. } t_{n-1, 1-\frac{\alpha}{2}} \in [-|t_0|, |t_0|]. \quad \text{Let } T \sim t_{n-1, 1-\frac{\alpha}{2}}.$$

The p -value (probability of being wrong) $p = 1 - \mathbf{P}(T \in [-|t_0|, |t_0|]) = 2 \times \mathbf{P}(T > |t_0|)$.