

PREDICTION OF THE RESPONSE VARIABLE

We observe (x_i, y_i) for $i=1, \dots, n$.

Consider an additional unit observed at a value x_* . We want to make a prediction about the value of the response variable corresponding to x_* .

The model is $Y_i = \beta_1 + \beta_2 x_i + \epsilon_i$, i.e. $E[Y_i] = \mu_i = \beta_1 + \beta_2 x_i$

hence $Y_* = \beta_1 + \beta_2 x_* + \epsilon_*$ with $\mu_* = \beta_1 + \beta_2 x_*$

The predicted value is $\hat{y}_* = \hat{\beta}_1 + \hat{\beta}_2 x_*$

The prediction \hat{y}_* corresponds to the estimate of the mean μ_* .

If we consider the estimators $\hat{\beta}_1$ and $\hat{\beta}_2$, we obtain the corresponding estimator $\hat{\mu}_*$ of the mean of Y_* .

We can study the distribution of $\hat{\mu}_*$.

$$\begin{aligned}\hat{\mu}_* &= \hat{\beta}_1 + \hat{\beta}_2 x_* = \bar{Y} - \hat{\beta}_2 \bar{x} + \hat{\beta}_2 x_* \\ &= \bar{Y} + \hat{\beta}_2 (x_* - \bar{x}) \\ &= \frac{1}{n} \sum_{i=1}^n Y_i + (x_* - \bar{x}) \sum_{i=1}^n w_i Y_i \quad \text{since } \hat{\beta}_2 = \frac{\sum_{i=1}^n w_i Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \text{ with } w_i = \frac{(x_i - \bar{x})}{\sum_{k=1}^n (x_k - \bar{x})^2} \text{ (see Lec. 2)} \\ &= \sum_{i=1}^n \left(\frac{1}{n} + (x_* - \bar{x}) w_i \right) Y_i\end{aligned}$$

$\Rightarrow \hat{\mu}_*$ is a linear combination of Y_1, \dots, Y_n

$\Rightarrow \hat{\mu}_*$ has normal distribution $\hat{\mu}_* \sim N(\dots, \dots) \rightarrow$ we need to find the mean and variance

$$E[\hat{\mu}_*] = E[\hat{\beta}_1 + \hat{\beta}_2 x_*] \stackrel{\text{linearity}}{=} \beta_1 + \beta_2 x_* = \mu_* \quad \text{unbiased}$$

$$\begin{aligned}\text{var}(\hat{\mu}_*) &= \text{var}\left(\sum_{i=1}^n \left(\frac{1}{n} + (x_* - \bar{x}) w_i\right) Y_i\right) \stackrel{\text{homoscedastic var}(Y_i) = \sigma^2}{=} \sum_{i=1}^n \left(\frac{1}{n} + (x_* - \bar{x}) w_i\right)^2 \sigma^2 = \\ &= \sum_{i=1}^n \left(\frac{1}{n^2} + w_i^2 (x_* - \bar{x})^2 + \frac{2}{n} w_i (x_* - \bar{x})\right) \sigma^2 = \\ &= \frac{1}{n} \sigma^2 + \sigma^2 (x_* - \bar{x})^2 \sum_{i=1}^n w_i^2 + 2 \sigma^2 (x_* - \bar{x}) \sum_{i=1}^n w_i = \\ &= \sigma^2 \left(\frac{1}{n} + (x_* - \bar{x})^2 \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) = \\ &= \sigma^2 \left(\frac{1}{n} + \frac{(x_* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)\end{aligned}$$

$$\Rightarrow \hat{\mu}_* \sim N\left(\mu_*; \underbrace{\sigma^2 \left(\frac{1}{n} + \frac{(x_* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}_{V(\hat{\mu}_*)}\right) = N(\mu_*, V(\hat{\mu}_*))$$

Let's derive a confidence interval for μ_*

We need a pivotal quantity

$$\frac{\hat{\mu}_* - \mu_*}{\sqrt{V(\hat{\mu}_*)}} \sim N(0, 1)$$

since $V(\hat{\mu}_*)$ involves the unknown σ^2 , similarly to what we have done for $\hat{\beta}_j$ we substitute $V(\hat{\mu}_*)$ with $\hat{V}(\hat{\mu}_*)$, obtaining

$$\frac{\hat{\mu}_* - \mu_*}{\sqrt{\hat{V}(\hat{\mu}_*)}} \sim t_{n-2} \quad \text{where} \quad \hat{V}(\hat{\mu}_*) = S^2 \left(\frac{1}{n} + \frac{(x_* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$$

Thus, a confidence interval of level $1-\alpha$ for μ_* is obtained as

$$1-\alpha = P\left(-t_{n-2; 1-\frac{\alpha}{2}} < \frac{\hat{\mu}_* - \mu_*}{\sqrt{\hat{V}(\hat{\mu}_*)}} < t_{n-2; 1-\frac{\alpha}{2}}\right)$$

$$1-\alpha = P\left(\hat{\mu}_* - t_{n-2; 1-\frac{\alpha}{2}} \sqrt{\hat{V}(\hat{\mu}_*)} < \mu_* < \hat{\mu}_* + t_{n-2; 1-\frac{\alpha}{2}} \sqrt{\hat{V}(\hat{\mu}_*)}\right)$$

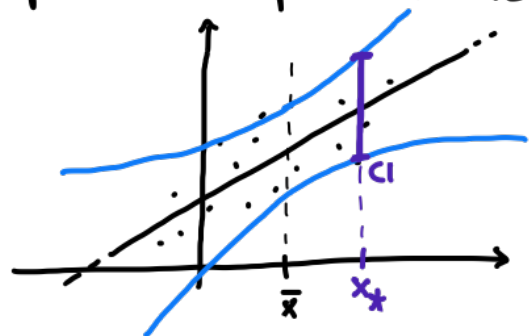
$$1-\alpha = P\left(\hat{\mu}_* - t_{n-2; 1-\frac{\alpha}{2}} \sqrt{S^2 \left(\frac{1}{n} + \frac{(x_* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)} < \mu_* < \hat{\mu}_* + t_{n-2; 1-\frac{\alpha}{2}} \sqrt{S^2 \left(\frac{1}{n} + \frac{(x_* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}\right)$$

conditioning now to the observed data: \hat{y}_* estimate of μ_* , S^2 estimate of σ^2

$$CI: \hat{y}_* \pm t_{n-2; 1-\frac{\alpha}{2}} \sqrt{S^2 \left(\frac{1}{n} + \frac{(x_* - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}$$

notice that the further x_* is from \bar{x} , the larger the CI will get

If I compute several pointwise CIs for varying x_* , I obtain "confidence bands"



(careful: the level $(1-\alpha)$ only holds pointwise)

\rightarrow why predicting outside of the range of the x_i is dangerous

These methods can be useful to formalize practical questions, for example:

- what is a reasonable set of values for Y if $x = \bar{x}$? \rightarrow compute CI for $\tilde{\mu}$
- is μ_0 a reasonable value for Y if I observe $x = \bar{x}$? \rightarrow test $H_0: \tilde{\mu} = \mu_0$
 $H_1: \tilde{\mu} \neq \mu_0$