

Simple Linear Regression (Ch 2)
Week 4 – Thursday
Applied Regression Analysis (STAT 757)

Paul J. Hurtado

Thursday, 11 Feb, 2016

Recap: Confidence Intervals

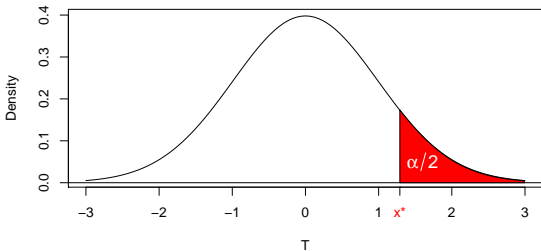
Intercept CI: The distribution of the intercept estimator $\hat{\beta}_0$ can be computed by *un-standardizing* the following r.v., which follows a *Student's t* distribution with $n - 1$ d.f.

$$T = \frac{\hat{\beta}_0 - \beta_0}{\text{se}(\hat{\beta}_0)}$$

Recall $S = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$ and $S_{XX} = \sum_{i=1}^n (x_i - \bar{x})^2$, and that the standard error of $\hat{\beta}_0$ is $\text{se}(\hat{\beta}_0) = S \sqrt{\frac{1}{n} \frac{\bar{x}^2}{S_{XX}}}$.

CI: 100(1- α)% of the time, *parameter* β_0 is in
 $\left[\hat{\beta}_0 + qt\left(\frac{\alpha}{2}, n - 2\right)\text{se}(\hat{\beta}_0), \hat{\beta}_0 + qt\left(1 - \frac{\alpha}{2}, n - 2\right)\text{se}(\hat{\beta}_0) \right]$

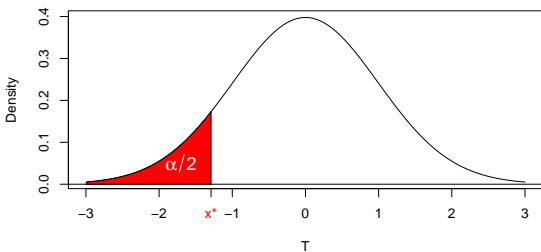
Note: $t()$ vs $qt()$



The T -score (and Z -score) tables found in classical textbooks tell you the *random variable value* x_* that satisfy

$$P(X > x_*) = \alpha/2$$

Our textbook uses $x_* = t(\alpha/2, n - 2)$ to denote these “upper tail” values. **However, in R...**



... we compute these values as **quantiles**. For example, the 25% quantile, x_* , satisfies $P(X \leq x_*) = 0.25$.

Thus, to clarify **textbook** vs **R** notation:

$$\begin{aligned} -t(\alpha/2, n - 2) &= \text{qt}(\alpha/2, n - 2) \\ t(\alpha/2, n - 2) &= \text{qt}(1 - \alpha/2, n - 2) \end{aligned}$$

Recap: Confidence Intervals

Slope CI: The distribution of the slope estimator $\hat{\beta}_1$ can be computed by *un-standardizing* the t_{n-1} distributed r.v.

$$T = \frac{\hat{\beta}_1 - \beta_1}{\text{se}(\hat{\beta}_1)}$$

Recall $S = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$ and $S_{XX} = \sum_{i=1}^n (x_i - \bar{x})^2$, and that the standard error of $\hat{\beta}_1$ is $\text{se}(\hat{\beta}_1) = S/\sqrt{S_{XX}}$.

CI: 100(1- α)% of the time, *parameter* β_1 is in
 $\left[\hat{\beta}_1 + qt\left(\frac{\alpha}{2}, n-2\right)\text{se}(\hat{\beta}_1), \hat{\beta}_1 + qt\left(1 - \frac{\alpha}{2}, n-2\right)\text{se}(\hat{\beta}_1) \right]$

Recap: Confidence Intervals

Regression Line CI: The distribution of $\hat{y}_* = \hat{\beta}_0 + \hat{\beta}_1 x_*$ (or, $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$) can be computed by *un-standardizing* the t_{n-1} distributed r.v.

$$T = \frac{\hat{y}_* - (\beta_0 + \beta_1 x_*)}{S \sqrt{\frac{1}{n} + \frac{(x_* - \bar{x})^2}{S_{XX}}}}$$

Thus,

CI: $100(1-\alpha)\%$ of the time, $E(\hat{y}_*) = \beta_0 + \beta_1 x_*$ is in

$$\hat{y}_* \pm qt\left(1 - \frac{\alpha}{2}, n - 2\right) S \sqrt{\frac{1}{n} + \frac{(x_* - \bar{x})^2}{S_{XX}}}$$

Recap: Prediction Intervals

Prediction Interval: The distribution of $\widehat{Y}_* = \widehat{\beta}_0 + \widehat{\beta}_1 x_*$ (or, $\widehat{y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 x_i$) can be computed by *un-standardizing* the t_{n-1} distributed r.v.

$$T = \frac{\widehat{y}_* - (\beta_0 + \beta_1 x_*)}{S \sqrt{1 + \frac{1}{n} + \frac{(x_* - \bar{x})^2}{S_{XX}}}}$$

Note the CI for $E(\widehat{y}_*)$ uses $se(\widehat{y}_*)$ while the prediction interval uses $se(\widehat{Y}_* - y_*)$. Thus,

CI: $100(1-\alpha)\%$ of the time, $E(\widehat{y}_*) = \beta_0 + \beta_1 x_*$ is in

$$\widehat{y}_* \pm qt\left(1 - \frac{\alpha}{2}, n - 2\right) S \sqrt{1 + \frac{1}{n} + \frac{(x_* - \bar{x})^2}{S_{XX}}}$$

Analysis of Variance

Observe that, if $\beta_1 = 0$ then the SLR model becomes

$$Y = \beta_0 + \epsilon \sim \text{Normal}(\beta_0, \sigma)$$

and so $\hat{\beta}_0 = \hat{y} = \bar{y}$. To test for a significant linear relationship, we test against this null hypothesis ($H_0 : \beta_1 = 0$) using

$$T = \frac{\hat{\beta}_1 - 0}{\text{se}(\hat{\beta}_1)} \sim t_{n-2}$$

In multiple regression, however, we need to generalize. This leads us to a different test statistic...

Analysis of Variance

Q: How much of the variation in y_i values comes from the linear component?

$$SST = S_{YY} = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$SS_{reg} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

It can be shown (see next slide) that

$$\begin{array}{ccccc} \text{Total variation in Y} & & \text{SS explained by regression} & & \text{residuals} \\ \underbrace{SST} & = & \underbrace{SS_{reg}} & + & \underbrace{RSS} \end{array}$$

Proof that $SST = SS_{reg} + RSS$:

$$\begin{aligned} SST &= \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n \overbrace{(y_i - \hat{y}_i + \hat{y}_i - \bar{y})}^{e_i}{}^2 \\ &= \sum_{i=1}^n e_i^2 + (\hat{y}_i - \bar{y})^2 + 2 e_i (\hat{y}_i - \bar{y}) = SS_{reg} + RSS + 2 \sum_{i=1}^n e_i (\hat{y}_i - \bar{y}) \end{aligned}$$

However, from our derivation of $\hat{\beta}_0$ and $\hat{\beta}_1$ (see textbook pg. 18), recall that $\sum_{i=1}^n e_i = 0$ and $\sum_{i=1}^n x_i e_i = 0$. Thus we see that

$$\begin{aligned} \sum_{i=1}^n e_i (\hat{y}_i - \bar{y}) &= \sum_{i=1}^n \hat{y}_i e_i - \bar{y} e_i = \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_i) e_i - \bar{y} \sum_{i=1}^n e_i \\ &= \hat{\beta}_0 \sum_{i=1}^n e_i + \hat{\beta}_1 \sum_{i=1}^n x_i e_i - \bar{y} \sum_{i=1}^n e_i = 0. \end{aligned}$$

Analysis of Variance

Generalizing the t test above, we can also get a p -value for the hypothesis test

$$H_0 : \beta_1 = 0 \quad \text{vs} \quad H_A : \beta_1 \neq 0$$

by using a more general test statistic that quantifies how much variation in Y results from the linear trend relative to the random variation determined by the magnitude of σ :

$$F = \frac{SS_{reg}/1}{RSS/(n-2)}$$

F has an $F_{1,n-2}$ distribution, and will be revisited in Ch. 5.

0-1 Categorical SLR

§2.6-2.7 on Monday