

Math 4501 - Probability and Statistics II

7.4 - Sample size

Overview: confidence intervals for a proportion

We will see how to select the sample size in order to estimate an unknown proportion to a given degree of accuracy.

\uparrow
maximum error E .

Confidence intervals for a proportion

Problem: Denote by $\hat{p} = \frac{y}{n}$ the proportion of observed successes in n trials and choose the sample size n so that the $100(1 - \alpha)\%$ confidence interval for p

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

is not longer than $\bar{x} \pm \varepsilon$, for some fixed $\varepsilon > 0$.

$\bar{x} = \hat{p} = \frac{y}{n}$ ← max error

Based on Approach 1 to estimate p (sec. 7.1-3)



Partial solution: Given the maximum error ε , we would like to solve for n the inequality

$$\frac{z_{\alpha/2} \sqrt{\hat{p}(1 - \hat{p})}}{\sqrt{n}} \leq \varepsilon$$

solve for n ??

Difficulty: we do not know what $\hat{p} = \frac{y}{n}$ actually is until we perform the experiment, and we first need to decide what the sample size n needs to be.

First solution: If, based on past experience, we know of a reasonable approximate p^* for the unknown proportion p , then we can set $\hat{p} = p^*$ in the maximum error of the point estimate to get

$$\frac{z_{\alpha/2} \sqrt{p^* (1 - p^*)}}{\sqrt{n}} \leq \varepsilon .$$

) solve for n

$$n \geq \frac{z_{\alpha/2}^2 p^* (1 - p^*)}{\varepsilon^2} .$$

Solving for n we get

and then take n at least as large as the smallest integer greater than the quantity on the right hand side of the inequality above.

Second solution: Solve

for n to get

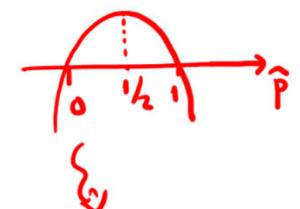
Note that

$$\hat{p}(1 - \hat{p}) \leq \frac{1}{4} \quad \text{worse possible case}$$

and use the maximum value of $1/4$ to get \Rightarrow choose n such that

$$n \geq \frac{z_{\alpha/2}^2}{4\varepsilon^2}.$$

$$\hat{p}(1 - \hat{p}) = \hat{p} - \hat{p}^2$$



Maximum occurs when $\hat{p} = 1/2$ and is equal to $\frac{1}{2}(1-\frac{1}{2}) = \frac{1}{4}$

Take n at least as large as the smallest integer greater than the quantity on the right hand side of the inequality above.

Note: this second approach is easier to implement, but yields a larger sample size than the first due to the lack of information regarding \hat{p} .

$$p^* = 0.08$$

Example

Suppose it is known that the unemployment rate on a certain country has been about 8%

An updated estimate is needed. How large should the sample be so that a 99% confidence interval provides an estimate for the actual unemployment rate to within 0.001?

$$\text{max error is } \varepsilon = 0.001$$

Let y be the number of unemployed people in a sample of size n. \leftarrow to be determined

Using that $z_{0.005} = 2.576$, we find that the interval is given by

$$z_{\alpha/2} \text{ with } \alpha = 0.01$$

$$\frac{y}{n} \pm 2.576 \sqrt{\frac{(y/n)(1 - y/n)}{n}}$$

According to Sec 7.3
endpoints of a 99%
confidence interval

The margin of error is

$$\hat{p} = \frac{y}{n} \approx p^* = 0.08$$

$$2.576 \sqrt{\frac{(y/n)(1 - y/n)}{n}} \approx 2.576 \sqrt{\frac{(0.08)(0.92)}{n}},$$

where the right hand side is obtained using the knowledge that the unemployment rate is close to 8%.

Solving

half width of interval using $\hat{p}^* = 0.08$

$$2.576 \sqrt{\frac{(0.08)(0.92)}{n}} \leq \underline{\underline{0.001}}$$

for n , we get

$$\boxed{n \geq 488,394},$$

which is a very large sample size...

To bring the sample to a more manageable size, we can:

- decrease the confidence level (from 99% to, for instance, 98%) $\rightarrow \alpha = 0.02$
- increase the allowed error (from 0.001 to, for instance, 0.01) $\rightarrow \varepsilon = 0.01$

Such choices would yield (note that $z_{0.01} = 2.236$):

$$z_{\alpha/2} \quad \rightsquigarrow \quad 2.236 \sqrt{\frac{(0.08)(0.92)}{n}} \leq \underline{\underline{0.01}}$$

) repeat process
with new value

and

$$\boxed{n \geq 3,982},$$

a much more reasonable sample size.

Example

Without access to any additional information, how large should a sample be so that a 95% confidence interval for a proportion provides an estimate of the true proportion to within 0.03?

$$\alpha = 0.05$$

$$\rightarrow \max_{\text{min}} \varepsilon = 0.03$$

Since we have no information about the approximate value of the true proportion, we require that

$$n \geq \frac{z_{\alpha/2}^2}{4\varepsilon^2}.$$

Setting $\alpha = 0.05$ (and so $z_{0.025} = 1.96$) and $\varepsilon = 0.03$, we get

$$n \geq \frac{(1.96)^2}{4(0.03)^2} \approx 1067.11$$

and choose to take a sample of size at least 1068.

endpoints for confidence interval

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$\downarrow \quad z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq \varepsilon$$

$$\Rightarrow n \geq \frac{z_{\alpha/2}^2 \hat{p}(1-\hat{p})}{\varepsilon^2}$$

Since \hat{p} is unknown $\hat{p}(1-\hat{p}) \leq \frac{1}{4}$

Take $n \geq \frac{z_{\alpha/2}^2}{4\varepsilon^2}$

Math 4501 - Probability and Statistics II

7.6 - Confidence intervals for regression parameters

(related with Sec. 7.1)

}

based mostly on finding confidence
intervals for the mean.

Simple linear regression: review

Problem: Given the data points $(\underline{x}_1, \underline{y}_1), (\underline{x}_2, \underline{y}_2), \dots, (\underline{x}_n, \underline{y}_n)$, estimate the parameters α and β of the linear model

$$E[Y|x] = \alpha + \beta(x - \bar{x})$$

Standing assumption: for each particular value of x , the value of Y differs from its mean by a random amount $\varepsilon \sim N(0, \sigma^2)$.

Consequence: For each $i = 1, 2, \dots, n$, we have

$$Y_i = \alpha + \beta(x_i - \bar{x}) + \varepsilon_i,$$

where $\alpha + \beta(x_i - \bar{x})$ are nonrandom and $\varepsilon_i \sim N(0, \sigma^2)$ are independent (with unknown variance σ^2).

The random variables $\underline{Y}_1, \underline{Y}_2, \dots, \underline{Y}_n$ are mutually independent normal variables

$$Y_i \sim N(\alpha + \beta(x_i - \bar{x}), \sigma^2), \quad i = 1, 2, \dots, n$$

Proposition (Maximum likelihood estimators of α , β , and σ^2) [From Chp. 6]

Under the conditions described above, the maximum likelihood estimators of α , β and σ^2 are given by:

$$\hat{\alpha} = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \rightarrow \text{linear combination of normally dist. r.v.}$$

$$\hat{\beta} = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n [Y_i - \hat{\alpha} - \hat{\beta}(x_i - \bar{x})]^2 = \frac{1}{n} \sum_{i=1}^n [Y_i - \hat{Y}_i]^2$$

where $\hat{Y}_i = \hat{\alpha} + \hat{\beta}(x_i - \bar{x})$

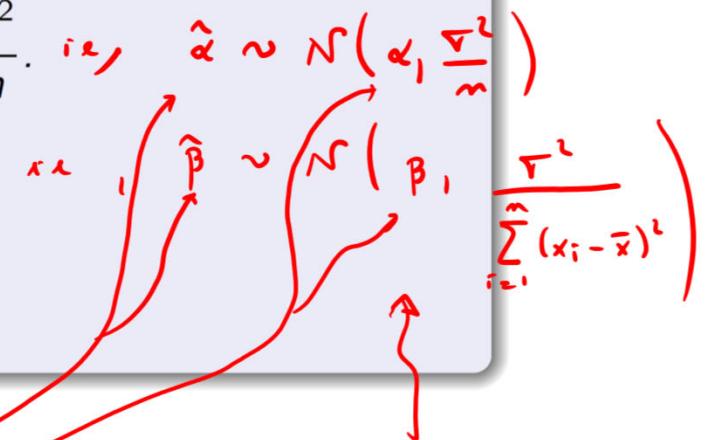
Proposition (Distributions of $\hat{\alpha}$ and $\hat{\beta}$) [From Chp. 5]

Under the conditions described earlier, we have that:

Proof
Chp. 6

- { 1) $\hat{\alpha}$ is normally distributed with mean $\underline{\alpha}$ and variance $\frac{\sigma^2}{n}$. i.e., $\hat{\alpha} \sim N(\alpha, \frac{\sigma^2}{n})$
- 2) $\hat{\beta}$ is normally distributed with mean β and variance $\frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$.

$$\text{Var}(\hat{\beta}) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$



Sec. 7.1

$$\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$$

$$\mu \in \bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}$$

Goal: to find
confidence intervals
for α and β !!!
BUT α and β are just
the unknown means of $\hat{\alpha}$ and $\hat{\beta}$

Distribution of $\hat{\sigma}^2$

Proposition [Proof is on the notes in Blackboard : very similar to the proof of $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$]

Under the conditions described earlier, we have that:

- 1) $\hat{\alpha}$, $\hat{\beta}$, and $\hat{\sigma}^2$ are mutually independent; \rightarrow analogue to \bar{x} and s^2 independent!
- 2) $\frac{n\hat{\sigma}^2}{\sigma^2}$ has a $\chi^2(n-2)$ distribution.

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \left[Y_i - (\hat{\alpha} + \hat{\beta}(x_i - \bar{x})) \right]^2$$

↑
↑
 $\hat{\alpha}$ and $\hat{\beta}$ are estimated parameters

compare with

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\frac{n\hat{\sigma}^2}{\sigma^2} = \sum_{i=1}^n \left[\frac{Y_i - (\hat{\alpha} + \hat{\beta}(x_i - \bar{x}))}{\sigma} \right]^2 \sim \chi^2(n-2)$$

$\frac{(n-1)S^2}{\sigma^2} = \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma} \right)^2 \sim \chi^2(n-1)$
-1 due to estimate \bar{x} for μ

Chp 5 { $x_i \sim N(\mu, \sigma^2) \Rightarrow \frac{x_i - \mu}{\sigma} \sim N(0, 1)$
 $\Rightarrow \left(\frac{x_i - \mu}{\sigma} \right)^2 \sim \chi^2(1) \Rightarrow \sum_{i=1}^n \left(\frac{x_i - \mu}{\sigma} \right)^2 \sim \chi^2(n)$

Confidence interval for α parameter of linear regression

We have seen that

✓ 1) $\hat{\alpha} \sim N\left(\alpha, \frac{\sigma^2}{n}\right)$ with σ^2 unknown

(α is the unknown mean from a normal distribution with unknown variance)

✓ 2) $\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi^2(n - 2)$

✓ 3) $\hat{\alpha}$ and $\hat{\sigma}^2$ are mutually independent

Then

and

and so

$$Z = \frac{\hat{\alpha} - \alpha}{\sqrt{\frac{\sigma^2}{n}}} \text{ is } N(0, 1)$$

$$U = \frac{n\hat{\sigma}^2}{\sigma^2} \text{ is } \chi^2(n - 2)$$

$$T_\alpha = \frac{Z}{\sqrt{U/(n-2)}} \text{ is } t(n-2).$$

$n-2$ is the # of d.o.f. of $U \sim \chi^2(n-2)$

Analogous
to Sec. 7.1

Note that

$$T_\alpha = \frac{Z}{\sqrt{U/(n-2)}} = \frac{\frac{\hat{\alpha} - \alpha}{\sqrt{\frac{\sigma^2}{n}}}}{\sqrt{\frac{n\hat{\sigma}^2}{\sigma^2(n-2)}}}$$

unknown σ^2 cancels out

$$T_\alpha = \frac{\hat{\alpha} - \alpha}{\sqrt{\frac{\hat{\sigma}^2}{n-2}}}.$$

can be rewritten as

Set

$$S_{\hat{\alpha}} = \sqrt{\frac{\hat{\sigma}^2}{n-2}}.$$

then

$$\boxed{T_\alpha = \frac{\hat{\alpha} - \alpha}{S_{\hat{\alpha}}} \sim t(n-2)}$$

Pick $t_0 = t_{\gamma/2}(n - 2)$ and note that

$$P \left(-t_0 \leq \frac{\hat{\alpha} - \alpha}{S_{\hat{\alpha}}} \leq t_0 \right) = 1 - \gamma$$

) solve for α

can be rewritten as

$$P(\hat{\alpha} - t_0 S_{\hat{\alpha}} \leq \alpha \leq \hat{\alpha} + t_0 S_{\hat{\alpha}})$$

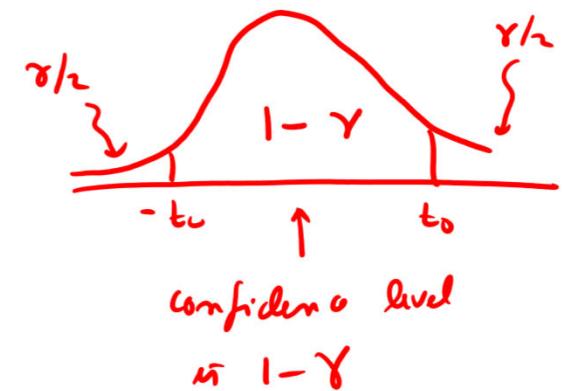
We conclude that

$$[\hat{\alpha} - t_0 S_{\hat{\alpha}}, \hat{\alpha} + t_0 S_{\hat{\alpha}}]$$

where

$$S_{\hat{\alpha}} = \sqrt{\frac{\hat{\sigma}^2}{n-2}}$$

is a $100(1 - \gamma)\%$ confidence interval for α .



Confidence interval for β or the other parameters of the linear regression

We have seen that

✓ 1) $\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)$ with σ^2 unknown

✓ 2) $\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi^2(n-2)$

✓ 3) $\hat{\beta}$ and $\hat{\sigma}^2$ are mutually independent

[we want a confidence interval for the mean β of a normal distribution with unknown variance]

suggest the use of t dist.

Then

$$Z = \frac{\hat{\beta} - \beta}{\sqrt{\frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \text{ is } N(0, 1)$$

and

$$U = \frac{n\hat{\sigma}^2}{\sigma^2} \text{ is } \chi^2(n-2)$$

and so

$$T_\beta = \frac{Z}{\sqrt{U/(n-2)}} \text{ is } t(n-2).$$

of dof of U

Note that

$$T_\beta = \frac{Z}{\sqrt{U/(n-2)}} = \frac{\frac{\hat{\beta} - \beta}{\sigma^2}}{\sqrt{\frac{n\sigma^2}{(n-2)}}}$$

can be rewritten as

$$T_\beta = \frac{\hat{\beta} - \beta}{\sqrt{\frac{n\sigma^2}{(n-2)\sum_{i=1}^n(x_i - \bar{x})^2}}}.$$

Set

$$S_{\hat{\beta}} = \sqrt{\frac{n\sigma^2}{(n-2)\sum_{i=1}^n(x_i - \bar{x})^2}}$$

← depends only on known quantities

so that $T_\beta = \frac{\hat{\beta} - \beta}{S_{\hat{\beta}}} \sim t(n-2)$

Pick $t_0 = t_{\gamma/2}(n - 2)$ and note that

t_0

$$P \left(\frac{\hat{\beta} - \beta}{S_{\hat{\beta}}} \leq \underline{t_0} \right) = 1 - \gamma$$

solve for β

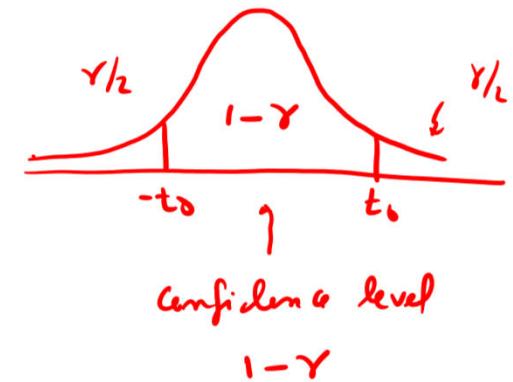
can be rewritten as

$$P \left(\underline{\hat{\beta} - t_0 S_{\hat{\beta}}} \leq \beta \leq \underline{\hat{\beta} + t_0 S_{\hat{\beta}}} \right)$$

We conclude that

$$\boxed{[\hat{\beta} - t_0 S_{\hat{\beta}}, \hat{\beta} + t_0 S_{\hat{\beta}}]}$$

is a $100(1 - \gamma)\%$ confidence interval for β .



when $S_{\hat{\beta}} = \sqrt{\frac{m \hat{\sigma}^2}{(n-2) \sum_{i=1}^n (x_i - \bar{x})^2}}$

Confidence interval for $\mu(x) = \alpha + \beta(x - \bar{x})$

Since

$$\hat{Y} = \hat{\alpha} + \hat{\beta}(x - \bar{x}) \quad \mu(x) = E[Y|x]$$

think of
 \hat{Y} as the
 point estimate
 for $\mu(x) = E[Y|x]$
 just like
 \bar{x} is for μ
 or $\hat{\beta}$ is for β

is a linear combination of the independent and normally distributed random variables $\hat{\alpha}$ and $\hat{\beta}$, \hat{Y} is also normally distributed.

The mean of \hat{Y} is

$$E[\hat{Y}] = E[\hat{\alpha} + \hat{\beta}(x - \bar{x})] = \underbrace{\hat{\alpha} + \beta(x - \bar{x})}_{\text{out of } \hat{Y}} + (x - \bar{x}) E[\hat{\beta}] = \mu(x)$$

and its variance is

$$\begin{aligned} \text{Var}(\hat{Y}) &= \text{Var}[\hat{\alpha} + \hat{\beta}(x - \bar{x})] = \text{Var}(\hat{\alpha}) + (x - \bar{x})^2 \text{Var}(\hat{\beta}) \\ &= \frac{\sigma^2}{n} + \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} (x - \bar{x})^2 \\ &= \sigma^2 \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \end{aligned}$$

$\hat{Y} \sim N(\mu(x), \sigma^2 \left(\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right))$

σ^2 is unknown!!

Since

$$Z = \frac{\hat{\gamma} - E[\hat{\gamma}]}{\sigma \sqrt{\frac{1}{n} + \frac{(x-\bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim N(0, 1)$$

$\sqrt{\text{Var}(\hat{\gamma})}$

$$\rightsquigarrow U = \frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi^2(n-2) ,$$

and independence $\hat{\alpha}$, $\hat{\beta}$, and $\hat{\sigma}^2$, allow us to conclude that

$$T_\mu = \frac{Z}{\sqrt{U/(n-2)}} = \frac{\frac{\hat{\alpha} + \hat{\beta}(x - \bar{x}) - [\alpha + \beta(x - \bar{x})]}{\sigma \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}}}{\sqrt{\frac{n\hat{\sigma}^2}{(n-2)\sigma^2}}} = \frac{\hat{\alpha} + \hat{\beta}(x - \bar{x}) - [\alpha + \beta(x - \bar{x})]}{\sqrt{\frac{n\hat{\sigma}^2}{n-2}} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim t(n-2)$$

has a $t(n-2)$ distribution.

Set $t_0 = t_{\gamma/2}(n - 2)$ and

$$S_\mu = \sqrt{\frac{n\widehat{\sigma}^2}{n-2}} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

so that $\bar{T}_\mu = \frac{\hat{y} - (\alpha + \beta(x - \bar{x}))}{S_\mu}$

and note that

$$P(-t_0 \leq \bar{T}_\mu \leq t_0) = 1 - \gamma$$

can be rewritten as

$$P(\underbrace{\hat{\alpha} + \hat{\beta}(x - \bar{x}) - t_0 S_\mu}_{\hat{y}} \leq \mu \leq \underbrace{\hat{\alpha} + \hat{\beta}(x - \bar{x}) + t_0 S_\mu}_{\hat{y}}) = 1 - \gamma$$

We conclude that

$$[\hat{\alpha} + \hat{\beta}(x - \bar{x}) - t_0 S_\mu, \hat{\alpha} + \hat{\beta}(x - \bar{x}) + t_0 S_\mu] \rightarrow \underline{\hat{y} \pm t_0 S_\mu}$$

is a $100(1 - \gamma)\%$ confidence interval for $\mu(x) = \alpha + \beta(x - \bar{x})$.

Note: Since S_μ depends on x , so does the width of the confidence interval constructed above.

Prediction interval

Problem: Having used the data points $(\underline{x}_1, \underline{y}_1), (\underline{x}_2, \underline{y}_2), \dots, (\underline{x}_n, \underline{y}_n)$ to estimate α , β , and σ^2 , we now want to estimate the value of Y corresponding to a given a value of $x = \underline{x}_{n+1}$.

$$\mu(x) = E[Y|x] = \alpha + \beta(x - \bar{x}) \rightarrow \text{estimate is } \hat{Y} = \hat{\alpha} + \hat{\beta}(x - \bar{x})$$

A point estimate of the corresponding value of \underline{Y} is given by

$$\hat{y}_{n+1} = \hat{\alpha} + \hat{\beta}(\underline{x}_{n+1} - \bar{x}),$$

but this is just one possible value of the random variable

$$Y_{n+1} = \alpha + \beta(\underline{x}_{n+1} - \bar{x}) + \underline{\varepsilon}_{n+1}.$$

The construction of a *prediction interval* for \underline{Y}_{n+1} when $x = \underline{x}_{n+1}$ is similar to that of the *confidence interval* for the mean of Y when $x = \underline{x}_{n+1}$.

previous case -

estimate of Y at a distinct value
of $x = \underline{x}_{n+1}$
prediction for the value
of Y at a distinct value
for x

Since

$$\text{def } \cup \overbrace{Y_{n+1}}^{\text{def } \cup} = \alpha + \beta(x_{n+1} - \bar{x}) + \varepsilon_{n+1} \quad \text{with} \quad \overbrace{\varepsilon_{n+1} \sim N(0, \sigma^2)}^{} , \Rightarrow Y_{n+1} \sim N\left(\underbrace{\alpha + \beta(x_{n+1} - \bar{x})}_{E[Y_{n+1}]}, \underbrace{\sigma^2}_{\text{Var}(Y_{n+1})}\right)$$

is normally distributed with mean $\alpha + \beta(x_{n+1} - \bar{x})$ and variance σ^2 , then

$$W = Y_{n+1} - \hat{\alpha} - \hat{\beta}(x_{n+1} - \bar{x}) ,$$

being a linear combination of independent and normally distributed random variables, is also normally distributed.

$$W \sim N\left(\underbrace{E[W]}_{=0}, \text{Var}(W)\right)$$

The mean of W is

$$\begin{aligned} E(W) &= E[Y_{n+1} - \hat{\alpha} - \hat{\beta}(x_{n+1} - \bar{x})] \\ &\stackrel{\text{linearity of } E}{=} E[Y_{n+1}] - E[\hat{\alpha}] - E[\hat{\beta}(x_{n+1} - \bar{x})] \\ &= \underbrace{\alpha + \beta(x_{n+1} - \bar{x})}_{=} - \underbrace{\alpha - \beta(x_{n+1} - \bar{x})}_{=} = 0 . \end{aligned}$$

Since $\underline{Y_{n+1}}$, $\underline{\hat{\alpha}}$, and $\underline{\hat{\beta}}$ are independent, the variance of \underline{W} is

$$\begin{aligned}
 \text{Var}(W) &= \text{Var} \left(\underline{Y_{n+1}} - \hat{\alpha} - \hat{\beta} (\underline{x_{n+1}} - \bar{x}) \right) \\
 &= \text{Var} (Y_{n+1}) + \text{Var} (\hat{\alpha}) + \text{Var} \left(\hat{\beta} (\underline{x_{n+1}} - \bar{x}) \right) \\
 &= \sigma^2 + \frac{\sigma^2}{n} + \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} (x_{n+1} - \bar{x})^2 \quad \text{Var}(\hat{\beta}) \\
 &= \sigma^2 \left[1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right].
 \end{aligned}$$

As a consequence, we get that

$$Z = \frac{W - \mathbb{E}[w]}{\sigma \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} = \frac{\underline{Y_{n+1}} - \hat{\alpha} - \hat{\beta} (\underline{x_{n+1}} - \bar{x})}{\sigma \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim N(0, 1)$$

$Z = \frac{W - E[w]}{\sqrt{\text{Var}(w)}}$

Recalling that

$$U = \frac{n\widehat{\sigma}^2}{\sigma^2} \sim \chi^2(n - 2),$$

independence of $\underline{Y_{n+1}}$, $\underline{\widehat{\alpha}}$, $\underline{\widehat{\beta}}$, and $\underline{\widehat{\sigma}^2}$, allows us to conclude that

$$\begin{aligned} T_Y &= \frac{Z}{\sqrt{U/(n-2)}} = \frac{\frac{Y_{n+1} - \widehat{\alpha} - \widehat{\beta}(x_{n+1} - \bar{x})}{\sigma \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}}}{\sqrt{\frac{n\widehat{\sigma}^2}{(n-2)\sigma^2}}} \xrightarrow{\text{from previous slide}} \\ &= \frac{Y_{n+1} - \widehat{\alpha} - \widehat{\beta}(x_{n+1} - \bar{x})}{\sqrt{\frac{n\widehat{\sigma}^2}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim t(n-2) \end{aligned}$$

has a $t(n - 2)$ distribution.

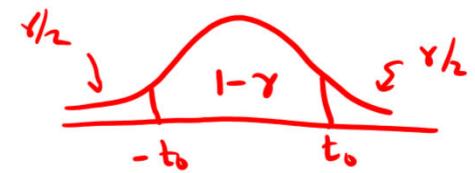
Set $t_0 = t_{\gamma/2}(n - 2)$ and

$$\text{so } S_Y = \sqrt{\frac{n\sigma^2}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

$$\text{so that } T_Y = \frac{Y_{n+1} - \hat{\alpha} - \hat{\beta}(x_{n+1} - \bar{x})}{S_Y}$$

and note that

$$P(-t_0 \leq T_Y \leq t_0) = 1 - \gamma$$



can be rewritten as

$$P(\hat{\alpha} + \hat{\beta}(x_{n+1} - \bar{x}) - t_0 S_Y \leq Y_{n+1} \leq \hat{\alpha} + \hat{\beta}(x_{n+1} - \bar{x}) + t_0 S_Y) = 1 - \gamma$$

We conclude that

$$[\hat{\alpha} + \hat{\beta}(x_{n+1} - \bar{x}) - t_0 S_Y, \hat{\alpha} + \hat{\beta}(x_{n+1} - \bar{x}) + t_0 S_Y]$$

$$\left. \begin{array}{c} \hat{\alpha} + \hat{\beta}(x_{n+1} - \bar{x}) \pm t_0 S_Y \\ \text{pt estimate for } Y_{n+1} \end{array} \right\}$$

is a $100(1 - \gamma)\%$ confidence interval for Y_{n+1} .