TMA4245 Statistikk

Week 15 - Wednesday

Week 14 (Mon. April 1 - Friday April 5)



If x has has linear relationship with y and x_i are fixed numbers:

1.
$$E[\varepsilon_1] = \ldots = E[\varepsilon_n] = 0 \Rightarrow$$

 $E[\hat{\alpha}] = \alpha$
 $E[\hat{\beta}] = \beta$

If x has has linear relationship with y and x_i are fixed numbers:

1.
$$E[\varepsilon_1] = \ldots = E[\varepsilon_n] = 0 \Rightarrow$$

 $E[\hat{\alpha}] = \alpha$
 $E[\hat{\beta}] = \beta$

2. $\varepsilon_1, \ldots, \varepsilon_n$ are iid with mean 0 and variance $\sigma^2 \Rightarrow$

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

If x has has linear relationship with y **3.** $\varepsilon_1, \ldots, \varepsilon_n \stackrel{iid}{\sim} N(0, \sigma^2) \Rightarrow$ and x_i are fixed numbers:

1.
$$E[\varepsilon_1] = \ldots = E[\varepsilon_n] = 0 \Rightarrow$$

 $E[\hat{\alpha}] = \alpha$
 $E[\hat{\beta}] = \beta$

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

2. $\varepsilon_1, \ldots, \varepsilon_n$ are iid with mean 0 and variance $\sigma^2 \Rightarrow$

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

If x has has linear relationship with y_{3} .. . and x_i are fixed numbers:

1.
$$E[\varepsilon_1] = \ldots = E[\varepsilon_n] = 0 \Rightarrow$$

 $E[\hat{\alpha}] = \alpha$

$$E[\hat{\beta}] = \beta$$

2. $\varepsilon_1, \ldots, \varepsilon_n$ are iid with mean 0 and variance $\sigma^2 \Rightarrow$

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

$$\varepsilon_{1}, \dots, \varepsilon_{n} \stackrel{n \sim n}{\sim} N(0, \sigma^{2}) \Rightarrow$$
$$\hat{\alpha} \sim N\left(\alpha, \frac{\sigma^{2} \sum_{i=1}^{n} x_{i}^{2}}{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}}\right)$$
$$\hat{\beta} \sim N\left(\beta, \frac{\sigma^{2}}{n \sum_{i=1}^{n} (x_{i} - \bar{x})^{2}}\right)$$

Notes:

- \triangleright $\hat{\alpha}$ and $\hat{\beta}$ are still approximately normal under **2.** provided *n* is large enough.
- In this class we assume 3. is necessary for simple linear regression unless otherwise stated

Assuming $\varepsilon_1, \ldots, \varepsilon_n \stackrel{iid}{\sim} N(0, \sigma^2)$:

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

Assuming $\varepsilon_1, \ldots, \varepsilon_n \stackrel{iid}{\sim} N(0, \sigma^2)$:

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

But σ^2 is unknown! Must instead be estimated. MLE (SME) is:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - (\hat{\alpha} + \hat{\beta} x_i))^2,$$

Assuming $\varepsilon_1, \ldots, \varepsilon_n \stackrel{iid}{\sim} N(0, \sigma^2)$:

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

But σ^2 is unknown! Must instead be estimated. MLE (SME) is:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - (\hat{\alpha} + \hat{\beta} x_i))^2,$$

What can we say about $\hat{\alpha}$, $\hat{\beta}$ if σ^2 is unknown?

Recall:

$$\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-2} \quad \text{and} \quad \frac{(n-2)S^2}{\sigma^2} \sim \chi^2_{n-2}$$
for $S^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - (\hat{\alpha} + \hat{\beta}x_i))^2$.

Also:

$$rac{Z}{\sqrt{rac{V}{
u}}}\sim t_{
u}$$

for $Z \sim N(0,1)$ and $V \sim \chi^2_{
u}$ independent.

Also, S^2 is independent of $\hat{\alpha}$, $\hat{\beta}$!

Hence, we get the following pivotal quantities under H_0 :

$$\frac{\hat{\alpha} - \alpha_0}{\sqrt{S^2 \frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}}} \sim t_{n-2}$$
$$\frac{\hat{\beta} - \beta_0}{\sqrt{S^2 \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim t_{n-2}.$$

Similar logic shows, for prediction at x_0 given by $\hat{\mu}_0 = \hat{\alpha} + \hat{\beta}x_0$:

$$E[\hat{\mu}_0] = \alpha + \beta x_0 = \mu_0$$

Var $(\hat{\mu}_0) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right),$

and:

$$\frac{\hat{\mu}_0 - \mu_0}{\sqrt{S^2\left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)}} \sim t_{n-2}.$$

Constructing confidence intervals

We can then use the pivotal quantities to construct $(1 - \alpha) \times 100\%$ confidence intervals:

$$\begin{aligned} \frac{\hat{\theta} - \theta_0}{\sqrt{S^2 \cdot C}} \sim t_{n-2} \\ -t_{\alpha/2,n-2} < \frac{\hat{\theta} - \theta_0}{\sqrt{S^2 \cdot C}} < t_{\alpha/2,n-2} \\ -t_{\alpha/2,n-2}\sqrt{S^2 \cdot C} < \hat{\theta} - \theta_0 < t_{\alpha/2,n-2}\sqrt{S^2 \cdot C} \\ -\hat{\theta} - t_{\alpha/2,n-2}\sqrt{S^2 \cdot C} < -\theta_0 < -\hat{\theta} + t_{\alpha/2,n-2}\sqrt{S^2 \cdot C} \\ \hat{\theta} - t_{\alpha/2,n-2}\sqrt{S^2 \cdot C} < \theta_0 < \hat{\theta} + t_{\alpha/2,n-2}\sqrt{S^2 \cdot C}, \end{aligned}$$

for constant C. This gives us the tools we need for hypothesis testing and confidence intervals!



Assume the following linear relationship between annual runoff Y and annual precipitation x within a drainage basin,

$$Y = \beta_0 + \beta_1 x + \varepsilon, \tag{1}$$

where β_0 and β_1 are unknown constants and ε is normally distributed with expected value (mean) 0 and unknown variance σ^2 .



Assume that we now are interested in predicting future runoff for a new year μ_0 given annual precipitation $x = x_0$, from the model defined in (1), where (x_0, Y_0) is independent of $(x_1, Y_1), (x_2, Y_2), \ldots, (x_{25}, Y_{25})$. Assume that $\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$ is a reasonable point estimator for the expected runoff $\mu_{Y|x_0} = \beta_0 + \beta_1 x_0$ when annual precipitation is x_0 . Assume $\hat{\beta}_0 = -1364$, $\hat{\beta}_1 = 1.08$, and $S^2 = 156^2$.



Assume $\hat{eta}_0=-1364$, $\hat{eta}_1=1.08$, and $S^2=156^2.$ Questions:

1. Show that $\hat{\mu}_0$ is unbiased

2. Show that
$$Var(\hat{\mu}_0) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$$



Assume $\hat{\beta}_0 = -1364$, $\hat{\beta}_1 = 1.08$, and $S^2 = 156^2$. Questions:

- 1. Show that $\hat{\mu}_0$ is unbiased
- 2. Show that $Var(\hat{\mu}_0) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 \bar{x})^2}{\sum_{i=1}^n (x_i \bar{x})^2} \right)$
- 3. Assume the average yearly precipitation in a year, \bar{x} , is 3200 mm/yr. Predict runoff for an average year.
- 4. Give a 95% confidence interval for the average runoff in a year with average precipitation.

Confidence intervals for the trend line

Recall:

$$E[\hat{\mu}_0] = \alpha + \beta x_0 = \mu_0$$

Var $(\hat{\mu}_0) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right),$

and:

$$\frac{\hat{\mu}_0 - \mu_0}{\sqrt{S^2\left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)}} \sim t_{n-2},$$

SO:

$$\hat{\mu}_{0} - t_{\alpha/2,n-2}\sqrt{S^{2} \cdot C} < \mu_{0} < \hat{\mu}_{0} + t_{\alpha/2,n-2}\sqrt{S^{2} \cdot C}$$

Confidence intervals for the trend line

95% confidence interval



The confidence interval (CI) margin of error (MOE) is then:

$$t_{\alpha/2,n-2}\sqrt{S^2 \cdot C} = t_{\alpha/2,n-2}\sqrt{S^2\left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)}$$

Confidence intervals for the trend line

95% confidence interval



The confidence interval (CI) margin of error (MOE) is then:

$$egin{aligned} t_{lpha/2,n-2}\sqrt{S^2\cdot C} &= t_{lpha/2,n-2}\sqrt{S^2\left(rac{1}{n}+rac{(x_0-ar{x})^2}{\sum_{i=1}^n(x_i-ar{x})^2}
ight)} \ &= t_{lpha/2,n-2}S\sqrt{rac{1}{n}+rac{(x_0-ar{x})^2}{\sum_{i=1}^n(x_i-ar{x})^2}}. \end{aligned}$$

What if the ϵ_i are non-Gaussian? If $\varepsilon_1, \ldots, \varepsilon_n$ iid, $E[\varepsilon_i] = 0$, and $Var(\varepsilon_i) = \sigma^2$, then for large enough *n*:

What if the ϵ_i are non-Gaussian? If $\varepsilon_1, \ldots, \varepsilon_n$ iid, $E[\varepsilon_i] = 0$, and $Var(\varepsilon_i) = \sigma^2$, then for large enough n:

 $\blacktriangleright~\hat{\beta}_0$ and $\hat{\beta}_1$ are still approximately Gaussian,

What if the ϵ_i are non-Gaussian?

If $\varepsilon_1, \ldots, \varepsilon_n$ iid, $E[\varepsilon_i] = 0$, and $Var(\varepsilon_i) = \sigma^2$, then for large enough *n*:

▶ $\hat{\beta}_0$ and $\hat{\beta}_1$ are still approximately Gaussian,

• $\hat{\mu}_0$ is still approximately Gaussian,

What if the ϵ_i are non-Gaussian?

If $\varepsilon_1, \ldots, \varepsilon_n$ iid, $E[\varepsilon_i] = 0$, and $Var(\varepsilon_i) = \sigma^2$, then for large enough *n*:

• $\hat{\beta}_0$ and $\hat{\beta}_1$ are still approximately Gaussian,

• $\hat{\mu}_0$ is still approximately Gaussian,

• S^2 still converges to $Var(\varepsilon)$ in some sense,

What if the ϵ_i are non-Gaussian?

If $\varepsilon_1, \ldots, \varepsilon_n$ iid, $E[\varepsilon_i] = 0$, and $Var(\varepsilon_i) = \sigma^2$, then for large enough *n*:

• $\hat{\beta}_0$ and $\hat{\beta}_1$ are still approximately Gaussian,

• $\hat{\mu}_0$ is still approximately Gaussian,

- S^2 still converges to $Var(\varepsilon)$ in some sense,
- ▶ approximate (1 − α) × 100% confidence intervals can still be created via:

$$\hat{\theta} - Z_{\alpha/2}\sqrt{S^2 \cdot C} < \theta_0 < \hat{\theta} + Z_{\alpha/2}\sqrt{S^2 \cdot C}$$

only using $Z_{\alpha/2}$ instead of $t_{\alpha/2,n-2},$ and

What if the ϵ_i are non-Gaussian?

If $\varepsilon_1, \ldots, \varepsilon_n$ iid, $E[\varepsilon_i] = 0$, and $Var(\varepsilon_i) = \sigma^2$, then for large enough *n*:

• $\hat{\beta}_0$ and $\hat{\beta}_1$ are still approximately Gaussian,

• $\hat{\mu}_0$ is still approximately Gaussian,

- S^2 still converges to $Var(\varepsilon)$ in some sense,
- approximate $(1 \alpha) \times 100\%$ confidence intervals can still be created via:

$$\hat{\theta} - Z_{\alpha/2}\sqrt{S^2 \cdot C} < \theta_0 < \hat{\theta} + Z_{\alpha/2}\sqrt{S^2 \cdot C}$$

only using $Z_{\alpha/2}$ instead of $t_{\alpha/2,n-2}\text{,}$ and

> approximate hypothesis tests can be conducted via the pivot:

 $\frac{\hat{\theta} - \theta_0}{\sqrt{S^2 \cdot C}} \sim Z \quad (`\sim' meaning `approximately distributed as')$

Last winter, on Sundays, Alexander sold cups of hot chocolate by the ski tracks near his house. This winter he plans to have a similar business. Alexander experienced that the sales changed dramatically with the weather and skiing conditions. He made a condition index, x, where x = 1 means "bad conditions", x = 2means "good conditions", x = 3 means "very good conditions" and x = 4 means "excellent conditions".

For 20 Sundays, i = 1, ..., 20, he registered both the number of cups sold, denoted y_i , and the associated condition index, x_i . We will phrase the sales as a regression model taking condition as an explanatory variable:

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, \dots, 20,$$

where $\epsilon_1, \ldots, \epsilon_{20}$ are independent variables with expected value 0 and variance σ^2 , and β_0 and β_1 are fixed but unknown regression parameters.

$$\sum_{i=1}^{20} (x_i - \bar{x})^2 = 24.95$$

$$\frac{1}{18} \sum_{i=1}^{20} (y_i - \hat{\beta}_0 - \hat{\beta}_1)^2 = 5.65^2$$

$$\hat{\beta}_1 \approx 9.51.$$

Questions:



$$\sum_{i=1}^{20} (x_i - \bar{x})^2 = 24.95$$

$$\frac{1}{18} \sum_{i=1}^{20} (y_i - \hat{\beta}_0 - \hat{\beta}_1)^2 = 5.65^2$$

$$\hat{\beta}_1 \approx 9.51.$$

Questions:

1. Would it be reasonable to construct a 95% confidence interval for β_0 or β_1 ?



$$\sum_{i=1}^{20} (x_i - \bar{x})^2 = 24.95$$

$$\frac{1}{18} \sum_{i=1}^{20} (y_i - \hat{\beta}_0 - \hat{\beta}_1)^2 = 5.65^2$$

$$\hat{\beta}_1 \approx 9.51.$$

Questions:

- 1. Would it be reasonable to construct a 95% confidence interval for β_0 or β_1 ?
- 2. Alexander thinks he can improve skiing conditions by 1 unit by grooming the snow, but for it to be worth it he would need to be 95% confident he would sell at least 8 more cups of hot chocolate on average. Use a 95% approximate confidence test to decide on your recommendation assuming $\varepsilon_1, \ldots, \varepsilon_n$ are iid with mean 0 and variance σ^2 .



Reminders

► Team based learning on Monday April 15