

CA09 - STAT100

Fábio Rodrigues Pereira

Oppgave 1

(Fra eksamen høst 2018. Det betyr at det nå er din tur å skrive den samme rapporten som du nettopp har lest eksempelbesvarelser av. Prøv å gjøre dette uten å se på de andre besvarelsene underveis, og sammenlign din/deres rapport med de gode besvarelsene etterpå.)

De pågående klimaendringene er et til stadig tilbakevendende tema og gjenstand for bekymring, både for vanlige mennesker, politikere og forskere. Hvordan vil en økning i temperatur påvirke naturtyper og økosystem?

Korallrev inneholder noen av de mest divergente økosystemer på planeten, de er habitat og ly for mange marine organismer. Revene har også en viktig rolle i å beskytte kystlinjer ved å ha en dempende effekt på bølgekraft og tropiske stormer, derfor er det ønskelig å beskytte korallrev så mye som mulig.

Det er gjort en studie på hvordan en økning i vanntemperaturer vil påvirke vekst i korallrev. Flere år på rad samlet forskere inn data på vanntemperaturer og hvor mye koraller hadde vokst i mm per år i Rødehavet.

Data som følger:

(x_i) Temperatur ($^{\circ}C$) Vekst (mm/år) (y_i)

29,7	2,63
29,9	2,57
30	2,67
30,2	2,6
30,5	2,47
30,7	2,39
30,9	2,25
31,2	2,24
31,5	2,15
28,9	2,77
29,4	2,53

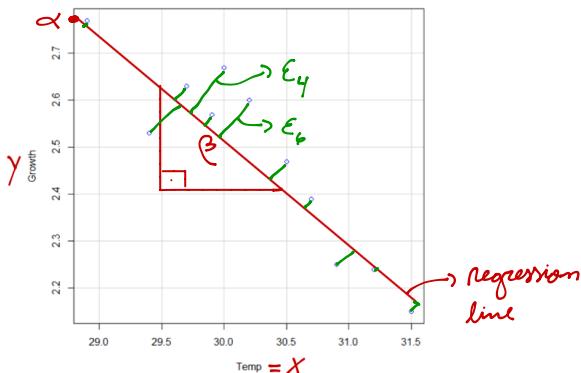
Det ble utført en lineær regresjonsanalyse. Nedenfor vises resultater fra SPSS:

	Coefficients:	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	9.50	1.06	9.35	0.0000062	
Temp	-0.23	0.03	-6.91	0.0000698	
	$\hat{\alpha}$	s	$\hat{\beta}$		
	s: 0.083 on 9 degrees of freedom				
	Multiple R-squared: 0.8415 $\rightarrow R^2$				
	Mean Temp		$\hat{X} = \bar{X} = \hat{\mu}_X$		
	(30.26)				

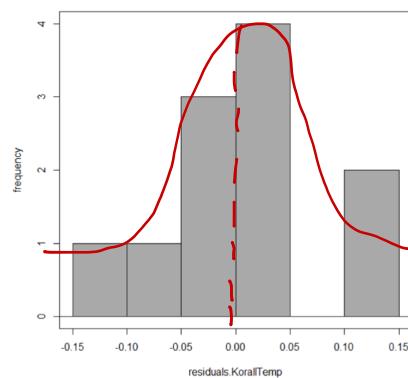
Skriv en kort rapport over analysen som er gjort (ta med modeller, antagelser, parametere med tolkning, estimer med usikkerhet, tester som kan være nyttige). Gi konkrete forklaringer på resultatene av forsøkene til en ikke-statistiker (for eksempel en klimaforsker).

→ It was performed a linear regression analysis where it is aimed to know how the increase of the global temperature can impact the ecosystems. The statistical model chosen is $Y_i = \alpha + \beta X_i + \epsilon_i$, where Y_i is the response variable, which represents the growth of the ocean reefs by mm/år, and correlated to the explanatory variable X_i which represents the independent yearly average temperature in $^{\circ}C$. α is the expected value of Y_i when the value of X_i is zero, also called the y-intercept of the regression line. β is another parameter which is the expected increase of Y_i when X_i increases 1 unit, also called slope of the regression line. ϵ_i is the residual or error of the model, for i observations between 1 to n , and which is independent normally distributed with mean 0 and standard deviation represented by σ . This first introduction can be technically illustrated as:

→ Model: $Y_i = \alpha + \beta X_i + \varepsilon_i$, where $\varepsilon_i \sim N(0, \sigma^2)$, independent for $i=1, \dots, n=11$, such that:



Figur 1: Årlig vekst av koraller mot vanntemperatur



Not completely symmetrical but near to normal distribution.

Now, we are interested to estimate the unknown parameters α , β and σ of the model. For that, we use the techniques called least squares to find unbiased estimates $\hat{\alpha}$ and $\hat{\beta}$, such that:

$$1.) \text{ Since } \bar{Y} = \frac{\sum_{i=1}^{11} Y_i}{11} \approx 2,48 \text{ and } \bar{X} = \frac{\sum_{i=1}^{11} X_i}{11} \approx 30,26, \text{ Then } \hat{\beta} = \frac{S_{xy}}{S_x^2} = \frac{\sum_{i=1}^{11} (Y_i - \bar{Y})(X_i - \bar{X})}{\sum_{i=1}^{11} (X_i - \bar{X})^2} \cdot \frac{(n-1)}{(n-1)} = -0,23 \text{ and } \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X} = 2,48 + 0,23 \cdot 30,26 \approx 9,50. \text{ we got } \hat{\alpha} = 9,5 \text{ and } \hat{\beta} = -0,23. \text{ It is still missing } \hat{\sigma}, \text{ which we will calculate below.}$$

2.) We know that the unbiased sample estimator S^2 for the estimate σ^2 is $S^2 = \frac{SSE}{n-2} = MSE$. Then, we need to split the variance of the model in order to find the SSE, such that:

↳ Since

$$\begin{cases} E[Y_i | X_i] = E[\alpha + \beta X_i + \varepsilon_i] = \alpha + \beta X_i + E[\varepsilon_i] \stackrel{*}{=} \alpha + \beta X_i = \hat{Y}_i \\ \text{Var}[Y_i | X_i] = \text{Var}[\alpha + \beta X_i + \varepsilon_i] = 0 + 0 + \sigma^2 = \hat{\sigma}^2 \\ \text{SD}[Y_i | X_i] = \sqrt{\hat{\sigma}^2} = \hat{\sigma} \end{cases}$$

and $Y_i = \hat{\alpha} + \hat{\beta} X_i + \varepsilon_i \rightarrow \varepsilon_i = Y_i - (\hat{\alpha} + \hat{\beta} X_i) = Y_i - \hat{Y}_i$

and $Y_i = \alpha + \beta X_i + \varepsilon_i \rightarrow Y_i - \bar{Y} = (\hat{\alpha} + \hat{\beta} X_i) - \bar{Y} + Y_i - (\hat{\alpha} + \hat{\beta} X_i) = \sum_{i=1}^{11} (Y_i - \bar{Y})^2 = \sum_{i=1}^{11} (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^{11} (\hat{Y}_i - Y_i)^2$

↳ Explained variance within the growth values of the coral reefs
↳ unexplained variance

② where SST is the total variance, SSR is the regression variance and SSE is the error/residual variance.

③ Then, the unbiased sample standard variation estimator S for the estimate σ is the square root of the mean of the error sum of squares, also called mean square error, such that:

$$S = \sqrt{S^2} = \sqrt{MSE} = \sqrt{\frac{SSE}{n-2}} = \sqrt{\frac{\sum_{i=1}^{11} [Y_i - (\hat{\alpha} + \hat{\beta} X_i)]^2}{9}} \approx 0,083 //$$

Now, we want to calculate how good is the model, where we use the coefficient of determination $R^2 = \frac{SSR}{SST} = 0,8415$. This score, which has range between $0 \leq R^2 \leq 1$, indicates that our regression has a very good precision, however we still need to perform the following tests to be secure of that.

We proceed a statistical test to check if the temperature has any effect to the coral growth, such that:

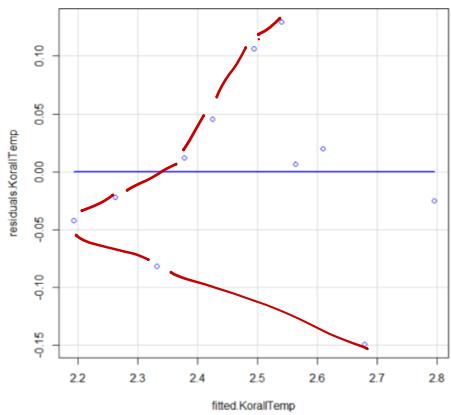
↳ Hypothesis claims: $H_0: \beta = 0$ where explanatory variable X (temperature) has no impact on Y ;
 $H_1: \beta \neq 0$ where " " " " " has impact on Y ;

↳ We choose a significance level $\alpha = 0,05$;

↳ Test statistic: $T^* = \frac{\hat{\beta} - \beta}{SE[\hat{\beta}]} \text{ where } SE[\hat{\beta}] = \sqrt{\frac{MSE}{\sum_{i=1}^n (x_i - \bar{x})^2}} \approx 0,03 \Rightarrow T^* = \frac{-0,23}{0,03} \approx -7,6 \Rightarrow |T^*| = 7,6 > 2,162 = t_{9,0,025}$
 ↳ Reject H_0

* Conclusion here: we are 95% secure that the temperature has impact on the coral reefs growth.

* Also, p-value $P(T \geq T^* | H_0 \text{ is true}) \approx 0,00007 < 0,05 = \alpha$ which also shows that the claim H_0 is rejected, in other words, we are more than 95% secure that the temperature has a negative effect ($\hat{\beta}$ is negative) on the growth of the coral reefs.



Figur 2. Residualer mot tilpassede verdier

→ Despite of our conclusion that the temperature has a linear combination which a regression model shows the negative impact of the temperature on the coral reefs growth, the diagnose plot of the residual values shows a problem. This plot should have gotten points near the 0 x-axis to confirm that the regression model is linear. However, its plot shows a expansion pattern which indicates that the model is not linear. This could have happened because of the scarce number of samples and would've been fixed by applying a transformation technique. □

En klimaforsker er bekymret for hvordan det vil gå med veksten i korallrev når temperaturen fortsetter å øke. Han ønsket å vite hva du anslår veksten av koraller vil være i mm per år når vanntemperaturen når 33°C . Hva sier du til han?

We can use $\hat{Y} = \hat{\alpha} + \hat{\beta} X_i$ to predict the value of the weight of the coral given $X_i = 33^\circ\text{C}$, such that:

$$\hookrightarrow \hat{Y} = 9,5 - (0,23 \cdot 33) \approx 1,90 \text{ mm per year.}$$

also, it is possible to give a prediction interval, with 95% of security, such that:

$$\hookrightarrow 95\% \text{ PI for } \hat{Y} = (\hat{\alpha} + \hat{\beta} X_i) \pm t_{df, \frac{\alpha}{2}} \cdot S \cdot \sqrt{1 + \frac{1}{m} + \left(\frac{X_i - \bar{X}}{S/SE[\hat{\beta}]} \right)^2}$$

$$\hookrightarrow = 1,90 \pm 2,262 \cdot 0,083 \cdot \sqrt{1 + \frac{1}{11} + \left(\frac{33 - 30,26}{0,083/0,03} \right)^2}$$

$$\hookrightarrow = [1,63; 2,17]$$

∴ We are 95% secure that the prediction interval will be between 1,63 to 2,17 mm/year given 33°C . □