


Regression - STAT110 Otago

Students also viewed

STAT110 STUDY TYPES


2 terms

 [Laura_Siegert57](#)

Preview

R CODES STAT110

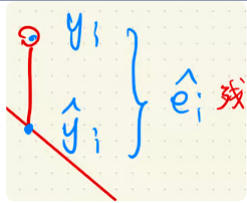
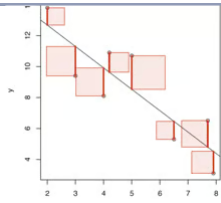
18 terms

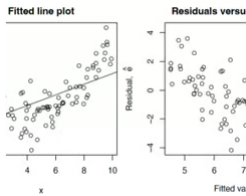
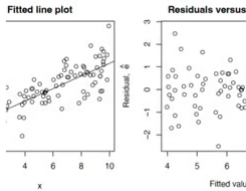
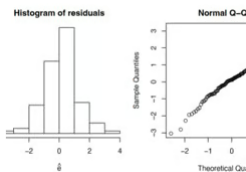
 [Laura_Siegert57](#)

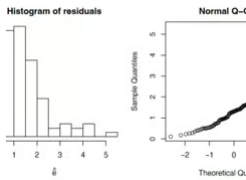
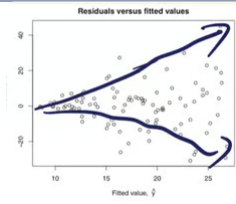
Preview

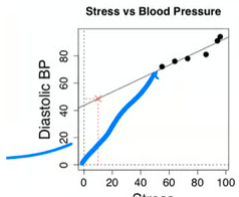
Terms in this set (60)

types of regression	linear (continuous data), logistic (categorical data), cox (categorical data in a survival analysis)
X	<i>Explanatory variable (X), also known as covariate, predictor, or independent variable.</i>
Y	<i>Outcome variable (Y), also known as response or dependent variable</i>
Simple linear regression (SLR)	<i>looks at a relationship between two continuous variables where the relationship between the two variables is approximately a straight line</i>

SLR equation	<p>$Y = \beta_0 + \beta_1 x + e$</p> <p>implies that the mean response is related to x by $\mu Y = \beta_0 + \beta_1 x$.</p> <p>$Y$ is the numerical outcome variable (continuous or approximately so)</p> <p>x is the explanatory variable</p> <p>β_0 is the intercept or constant (where the line crosses the y axis)</p> <p>β_1 is the slope of the line</p> <p>e (often denoted ϵ) is the random error or residual term</p>
SLR equation for estimating	$= \hat{\beta}_0 + \hat{\beta}_1 x$
residual ('estimated error') term	<p>$\hat{e}_i = y_i - \hat{y}_i$</p> 
How to find regression line	<p>The line of best fit minimises the sum of the squares of the residuals.</p> 
equation for how to find regression line	$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})}$
How to calculate beta1 and beta0	$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$ $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$

example for how to calculate regression (Stress and Blood Pressure)	<ol style="list-style-type: none"> 1. get n, n = 6 2. find the <i>explanatory and outcome</i> 3. calculate beta 1 and beta 0 4. get the regression equation 5. using R for SLR
Assumptions for Simple Linear Regression(LINE)	<p>Linearity: The relationship between the mean response μY and x is described by a straight line.</p> <p>Independence The responses Y_1, Y_2, \dots, Y_n are statistically independent.</p> <p>Normality The error terms e_1, e_2, \dots, e_n come from a normal distribution.</p> <p>Equal variance The errors terms all have the same variance, σ^2 ('homoscedastic').</p>
What diagram is used for checking linearity	<p>residual</p> $= y_i - \hat{y}_i$ $= y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$ $= y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$
failure of linearity assumption	
linearity assumption holds	
Checking independence assumption	<p>May get insight by thinking about the study design. (Ask yourself questions)</p>
plot for Checking the normality assumption	<p>Q - Q plot</p>
pass of normality assumption	

fail of <i>normality assumption</i>	
Checking <i>equal variance assumption (homoscedasticity)</i>	<p>pass if the residual plot is not like this</p> 
what is the impact if <i>Fail of the linearity assumption</i>	<p>critical</p> <p>If that assumption fails, all conclusions drawn from the model will be invalid.</p>
what is the impact if <i>Fail of independence or equal variance assumptions</i>	<p>remain valid</p> <p>However, estimates can be inefficient</p> <p>Follows that fitted regression line is useable.</p> <p>Any test results or confidence intervals based on the regression model will be invalid.</p>
what is the impact if <i>Fail of normality assumption</i>	<p>typically least important.</p> <p>Effects validity of confidence intervals and test results when the sample size n is small.</p>
what to do with outliers	<p>the first thing to do is check that the data are correctly recorded</p> <p>If data cannot be corrected, try refitting regression with outliers removed, but still investigate cause of outliers – may be very important.</p>
estimate of error variance	$s_e^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{e}_i^2$ <p>\hat{e}_i^2 is the residual</p>
(estimated) standard error	$\frac{s_e}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$
degree of freedom for SLR's CI	<p>$v = n - 2$</p> <p>because there're two parameters</p>

what is the multiplier for SLR's CI	$\frac{\text{estimate} - \text{std. err}}$
what is the SE for SLR's CI	$\sqrt{\sum_{i=1}^n}$
Using R to find SLR's CI	<pre>(model1) 2.5 % ;) 24.8300345 6 0.2557407</pre>
$\beta_1 = 0$ indicates what	<p>that the response is not (linearly) related to the predictor.</p> <p>so the estimated slope will (almost) always be non-zero: $\hat{\beta}_1 \neq 0$.</p>
Steps to test to assess strength of evidence in the data for $\beta_1 \neq 0$	<ol style="list-style-type: none"> 1. Setting up the hypotheses $H_0 : \beta_1 = 0, H_A : \beta_1 \neq 0$. 2. Calculating The test statistic (picture) 3. Computing the p-value 4. draw conclusion with rejecting or not H_0 $t = \frac{\hat{\beta}_1}{s_{\hat{\beta}_1}}.$
when predicting the data	ignore e_0
why not recommend to extrapolate when predicting data	<p>plot may not be linear</p> 
prediction error	<p>The prediction error is analogous to a standard error, but takes account of both sources of uncertainty.</p> <p>For prediction at x_0, the prediction error is:</p> $\sqrt{1 + \frac{1}{n} + \sum}$

prediction interval
formula

$$-\frac{\alpha}{2}, n-2)$$

correlation coefficient (r)

summarises the strength of a linear relationship between variables.

It is a measure of linear association between variables
It describes both the strength and direction of the relationship.

$$r \in [-1, 1]$$

A **positive value** of r means that **Y and X increase together**.

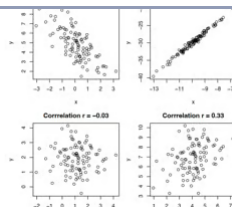
A **negative value** of r means that as **X increases, Y decreases** (and vice-versa).

The **strength of the linear relationship increases** as r tends **towards 1 or -1**.

r = 0 corresponds to no linear relationship between the variables.

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

scatterplots for r



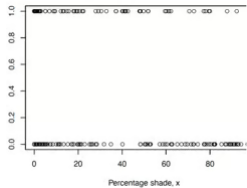
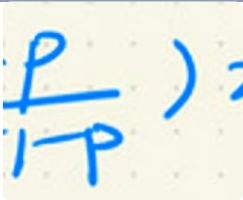
re-write for r

s_x and s_y are sample standard deviations for x and y variables.

s_{xy} is the sample covariance between x and y.

$$r = \frac{s_{xy}}{s_x s_y}$$

<p>Correlation coefficient versus regression models</p>	<p>The correlation coefficient is a summary of the data.</p> <p>Unlike linear regression, the correlation coefficient does not specify a model for the data, and cannot (for example) be used for prediction.</p> <p>The correlation coefficient is symmetric in the variables. That is, correlation between x and y is the same as correlation between y and x.</p> <p>In regression, the variables are not handled symmetrically.</p> <p>Regression models look at variation in Y for fixed values of x.</p>
<p>coefficient of determination (R^2)</p>	<p>R^2, is a measure of how well a regression model describes the data.</p> <p>R^2 is the squared correlation between the observed and predicted responses</p> <p>$R^2 \in [0, 1]$</p>
<p>meaning for the value of R^2</p>	<p>A high value of R^2 (close to 1) indicates a regression model that describes the data very well.</p> <p>Conversely, a low value of R^2 (close to 0) indicates a regression that describes the data poorly.</p> <div data-bbox="1114 1041 1364 1220"> </div>
<p>what describes the overall variation in the response variable?</p>	<p>total sum of squares</p> $= \sum_{i=1}^n (y_i - \bar{y})^2$
<p>what describes the total variation of the data points about the regression line?</p>	<p>residual sum of squares</p> <p>RSS can be thought of as variation not explained by the regression model</p> $= \sum_{i=1}^n (y_i - \hat{y}_i)^2$

what describes as the amount of variation in the response that is explained by the regression model?	<p>explained sum of squares</p> <p>ESS = TSS - RSS</p>
Equation of R ²	$\frac{ESS}{TSS} = 1 -$
Correlation does not equal causation	<p>e.g., just because there's more ice cream in the summer and more drowning in the summer doesn't mean there's a link between ice cream and drowning.</p>
logistic regression	<p>outcome variable is binary</p> 
equation for logistic regression	<p>Y is the binary outcome variable, $Y = 1$ or $Y = 0$ for each observation.</p> <p>p is the probability that specified category will occur; i.e. $p = \Pr(Y = 1)$.</p> <p>x is the explanatory variable.</p> <p>Parameters β_0, β_1 are the regression coefficients.</p> <p>β_0 is intercept and β_1 slope 'on the logit scale'</p> <p>In formula, log is the natural logarithm (log to base e).</p> $g\left(\frac{p}{1-p}\right)$
	
Which technique we use when we estimating the regression coefficients?	<p>maximum likelihood estimation</p>

what will <i>increasing x by one unit</i> results in?	a multiplicative change of e^{β_1} to the odds
formula for <i>logistic curve</i> for the probability p	$\frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$
Testing in logistic regression	<p>1. Define the hypotheses: $H_0 : \beta_1 = 0$ and $H_A : \beta_1 \neq 0$.</p> <p>2. The test statistic is: where $s_{\hat{\beta}_1}$ is the standard error of $\hat{\beta}_1$.</p> <p>The further away our test statistic, z, is from zero, the greater the evidence against H_0.</p> <p>3. get the corresponding p-value</p> <p>4. reject/not reject H_0</p> <p>5. conclusion</p> $z = \frac{\hat{\beta}_1}{s_{\hat{\beta}_1}}$
Multiple regression model	$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$
mean value de la Multiple regression model	$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$
Applications of multiple regression	<p>1 Adjusting for the effect of confounding variables.</p> <p>2 Establishing which variables are important in explaining the values of the response variable.</p> <p>3 Predicting values of the response variable.</p> <p>4 Describing the strength of the association between the response variable and the explanatory variables.</p>

least squares estimates	$\hat{\beta}_1(y_i -$
RSS for \hat{e}_i	<p>to estimate the error variance</p> $\sigma^2_e = \sum_i$
usual estimate	$= \frac{RSS}{n - k -$