

Topic: Inference for Regression

We have looked at inference for means and proportions, for one and two samples. We have looked at inference for categorical data. One of those tests was to see if there was an association between two categorical variables. In this unit, we look at inference procedures about the relationship between two *quantitative* variables measured on the same individuals.

Suppose we wanted to answer questions like:

“Is there really a relationship between arm-span and height in adults, or could the pattern we see in the scatterplot plausibly happen by chance?”

“In the population of adults, how much will the predicted height change for each increase of 1 inch in arm span? What is the margin of error for this estimate?”

Questions like these involve **regression**.

We will revisit the Normal and  $t$  distributions for our test statistic.

Procedures for inference still follow the State, Plan, Do, Conclude method. Some of the details are different, particularly the conditions for inference. We also learn some techniques to consider if some of the conditions are not met.

If you do not remember what we did with regression, then it would be a good idea to re-read Chapter 3.

Vocabulary you should recall from Chapter 3:

explanatory variable	Variable that may help explain or predict changes in a response variable. (Graphed on the horizontal axis.)
response variable	Variable that measures an outcome of a study. (Graphed on the vertical axis.)
scatterplot	Plot that shows the relationship between two quantitative variables measured on the same individuals.
correlation $r$	Measures the direction and strength of the linear relationship between two quantitative variables. (On formula sheet.)
direction	the direction of the association of two quantitative variables: positive, negative or something else (describe)
form	the form of the association of two quantitative variables: linear, curved, or something else (describe)
strength	the strength of the association of two quantitative variables: described such as: weak, moderate, strong or maybe none
outlier	Observation that lies outside the overall pattern of the other observations.
positive association	When above-average values of one variable to accompany above-average values of the other and also of below-average values to occur together.
negative association	When above-average values of one variable tend to accompany below-average values of the other.
regression line	Line that describes how a response variable $y$ changes as an explanatory variable $x$ changes.
predicted value	$\hat{y}$ (read “y hat”) is the predicted value of the response variable $y$ for a given value of the explanatory variable $x$ .
slope	In a regression line, the amount by which $y$ is predicted to change when $x$ increases by one unit.
y-intercept	In a regression line, the predicted value of $y$ when $x = 0$ .
extrapolation	Use of a regression line for prediction far outside the interval of values of the explanatory variable $x$ used to obtain the line. Such predictions are often not accurate.
residual	Difference between an observed value of the response variable and the value predicted by the regression line: $y - \hat{y}$
least-squares regression line (LSRL)	The line that makes the sum of the squared vertical distances of the data points from the line as small as possible.
residual plot	Scatterplot of the residuals against the explanatory variable. Residual plots help us assess whether a linear model is appropriate.
standard deviation of the residuals ( $s$ )	This value gives the approximate size of a “typical” prediction error (residual).
coefficient of determination: $r^2$	Fraction of the variation in the values of $y$ that is accounted for by the least-squares regression line of $y$ on $x$ .
influential observations	An observation is influential for a statistical calculation if removing it would markedly change the result of the calculation.

New Vocabulary introduced in Chapter 12:

$\mu_y = \alpha + \beta x$	True regression line based on the entire population of data.
$\sigma_b = \frac{\sigma}{\sigma_x \sqrt{n}}$	The standard deviation of the sampling distribution of $b$ .
$SE_b = \frac{s}{s_x \sqrt{n-1}}$	The standard error of the sampling distribution of $b$ . This approximates $\sigma_b$ using sample data.
L-I-N-E-R	<ul style="list-style-type: none"> <li>• Linear – the linear model <math>\mu_y = \alpha + \beta x</math> exists</li> <li>• Independent – the observations are independent</li> <li>• Normal – all the <math>y</math> distributions are Normal for each <math>x</math> value</li> <li>• Equal SD – all the <math>y</math> distributions have the same <math>\sigma</math>.</li> <li>• Random – the data come from a random sample or a randomized experiment</li> </ul>
$t$ interval for slope	Confidence interval for $\beta$ , the true slope of the regression line
$t$ test for slope	Significance test for $\beta$ , the true slope of the regression line.
$\rho$ (rho)	The true correlation $r$ measuring the direction and strength of the linear relationship between the two quantitative variables.
transforming data	Transforming data can sometimes change curved relationships between two quantitative variables into linear relationships.
power model	A curved relation of the form $y = ax^P$ . Can be transformed to achieve linearity by any of: $(x^P, y)$ , $(x, \sqrt[P]{y})$ , $(\ln x, \ln y)$ .
logarithm	In stat, we generally use natural logarithms: $e^A = B \Leftrightarrow \ln B = A$
exponential model	A curved relation of the form $y = ab^x$ . Can be transformed to achieve linearity by $(x, \ln y)$ .

Summary of Inference for  $\beta$ , True Slope of the Regression Line

Test of Significance	
State	$H_0: \beta = \beta_0$ $H_a: \beta < \beta_0$ or $H_a: \beta \neq \beta_0$ or $H_a: \beta > \beta_0$ <p>Where <math>\beta</math> is the true slope of the regression line of &lt;response variable&gt; on &lt;explanatory variable&gt;.</p> <p><math>\alpha = \dots</math> (select a significance level)</p>
Plan	<p>If conditions for inference are met, conduct a <math>t</math> test for linear regression.</p> <p>Linear: Check to see that a linear model fits the data; linear is a good fit if the scatterplot of the residuals show a patternless scatter.</p> <p>Independent: Check that the observations are independent. When sampling without replacement, check the 10% condition.</p> <p>Normal: Plot the residuals to make sure there is no strong skew or outliers; a probability plot should be approximately linear.</p> <p>Equal SD: Scatterplot of the residuals vs. the explanatory variable should show about the same (vertical) spread for every level of <math>x</math>.</p> <p>Random: Check that the data come from a random sample or a randomized experiment.</p>
Do	<p>the test statistic is <math>t</math> with <math>n - 2</math> degrees of freedom:</p> $t = \frac{b - \beta_0}{SE_b} \sim t(n - 2)$ <p>The P-value is <math>P(T &lt; t)</math> or <math>2P(T &gt;  t )</math> or <math>P(T &gt; t)</math> depending on the direction of your alternative hypothesis</p>
Conclude	<p>If P-Value is less than alpha (significance level), then we have evidence in favor of the alternative hypothesis.</p> <p>Otherwise, we do NOT have evidence in favor of the alternative hypothesis.</p>
Confidence Interval	<p>If conditions for inference are met, a CI% <math>t</math> interval for <math>\beta</math>, the true slope of the regression line of &lt;response variable&gt; on &lt;explanatory variable&gt; is given by:</p> $b \pm t^* SE_b$ <p>where the critical value <math>t^*</math> comes from a <math>t</math> distribution with <math>n - 2</math> degrees of freedom.</p>