

Chi-square test.

This test looks at the variables on the side and top axes of a table and tests whether they are independent.

For example it can be used to show whether variations in political opinions depend on a respondent's age.

H_0 No association between two variables

H_1 There is an association

Expected values: $\frac{\text{Row total} \times \text{column total}}{\text{total no. of observations}}$

χ^2 formula: $\sum_{\text{all cells}} \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$

$DF = (\text{number of rows} - 1) \times (\text{number of columns} - 1)$

Create confidence interval with DF.

Example: Sample paper.

If p value (eg 0.04) is less than significance level (0.05) we can reject null hypothesis that there is an association between two variables yes/no and male/female.

Using df create a confidence interval eg.

$$-3.84 < x < 3.84$$

If chi square value is less than 3.84, not enough evidence against H_0 .

95% interval for difference in proportion

$$= (p_{dec} - p_{jue}) \pm t_{critical} * SE(p_{dec} - p_{jue})$$

$$df = n_1 + n_2 - 2 = 196$$

$$\pm 1.96 * SE$$

p_1 and p_2 are the sample proportions $\frac{count}{total}$

$$\frac{\text{yes Decade}}{\text{total decade}} - \frac{\text{yes June}}{\text{total June}} = 0.02$$

$$0.02 \pm SE = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

$$= -0.02 \text{ to } 0.06$$

neg end mean month she June 2020

Regression

Scatterplot:

- Consists of x-axis and y-axis and a series of dots
- Each dot represents one observation from a dataset
- Useful for determining the correlation between two variables
- Independent variable on x-axis
- Dependent variable of y-axis

Example:

Regression analysis: Yield gram (Y-axis) Temperature (Celsius)

Minitab output: Yield gram = $17.0 + 2.00 \text{ temp}(C)$

Predictor	Coef	SE Coef	T	P
constant	17.002	4.072	4.18	0.000
temp C	1.9957	0.05334	37.41	0.000

S = 4.01967

R-Sq = 98.4%

R-Sq (adj) = 98.3%

Yield = equation made using coef rounded to nearest whole no

17.0 = yield when $t=0$

2.0 = increase in Y per unit increase in X

Residual plot = Observed val - Predicted val

Actual value - value got using equation

SE coef - measure of variability of the slope from sample to sample

Rough 95% interval for $2.00 \pm 1.96 * 0.05$

every degree C extra (1.99, 2.01) increase yield by (1.99 and 2.01)

dependent variable
P-values: implicitly test hypothesis: - (constant)

H_0 Population Slope = 0

(dependent)

H_1 Population Slope $\neq 0$

Test Statistic $\frac{\text{Coeff} - 0}{\text{SE(Coeff)}} = \frac{2.00}{0.0534} = 37.4451$ given by T

$p < 0.001$ we reject H_0

P value in independent variable: (Temp C)

Test hypothesis H_0 Population intercept = 0

H_1 Population intercept $\neq 0$

$p < 0.001$ we reject H_0

R^2 :

- Measure of the fit of the model to the data

- 98.4% of the variance of "constant" accounted for.

S:

- For each observation we calculate the residual.

- $S = 4.0197$ is the standard deviation of the residual

Assumptions:

- Residuals should be normally distributed

- There should be no relationship between the residuals and the predicted value