

Contents

1 Reading 9: Correlation and Regressions	1
1.1 Sample covar and sample correlation coefficients	1
1.2 Limitations to correlations analysis	1
1.3 Hypothesis: determine if the population correlation coefficient is zero	1
1.4 Determine dependent/independent variables in a linear regression	1
1.5 Assumptions in linear regression and interpret regression coeff.	1

1 Reading 9: Correlation and Regressions

1.1 Sample covar and sample correlation coefficients

Sample covariance: $cov_{x,y} = \sum_i \frac{(X_i - \bar{X})(Y_i - \bar{Y})}{n-1}$

Sample correlation coeff: $r_{x,y} = \frac{cov_{x,y}}{s_x s_y}$, where s_x is the sample dev of X.

1.2 Limitations to correlations analysis

Outliers: The results will be affected by extreme data points.(outliers)

Spurious correlation: There might be some non-zero correlation coeff, but actually they have no correlation at all.

Nonlinear relationships: Correlation only describe the linear relations.

1.3 Hypothesis: determine if the population correlation coefficient is zero

Two-tailed hypothesis test:

$$H_0 : \rho = 0, H_a : \rho \neq 0$$

Assume that the two populations are **normally** distributed, then we can use t-test:

$$t = \frac{r\sqrt{n-2}}{1-r^2}$$

: Reject H_0 if $t > +t_{critical}$ or $t < -t_{critical}$. Here, r is the sample correlation. Remember, you need to check t-table to find the t-value.

1.4 Determine dependent/independent variables in a linear regression

Simple linear regression: Explain the variation in a dependent variable in terms of the variation in a single independent variable. **Independent variables** are called explanatory variable, the exogenous variable, or the predicting variable. **Dependent variable** is also called the explained variable, the endogenous variable, or the predicted variable.

1.5 Assumptions in linear regression and interpret regression coeff.

1. Assumptions of linear regression:

- Linear relationship must exist.
- The independent variable is uncorrelated with residuals.
- Expected Residual term is value. $E(\epsilon) = 0$
- variance of the residual term is const. $E(\epsilon_i^2) = \sigma_\epsilon^2$
- The residual term is independently distributed. $E(\epsilon_i \epsilon_j) = 0$
- The residual term is normally distributed.

2. Simple Linear Regression Model

- (a) Model: $Y_i = b_0 + b_1 X_i + \epsilon_i$, where $i = 1 \dots n$, and Y_i is the actual observed data.
- (b) The fitted line, the line of best fit : $\hat{Y} = \hat{b}_0 + \hat{b}_1 X_i$. Where \hat{b}_0 is the estimated parameter of the model.
- (c) How to choose the best fitted line? **Sum of squared errors** is minimum.

$$\hat{b}_1 = \frac{cov_{x,y}}{sigma_x^2}$$

where X is the indepdent variable.

$$\hat{b}_0 = \bar{Y} - \hat{b}_1 \bar{X}$$

where \bar{X}, \bar{Y} are the mean.