

1 Simple Linear Regression: Background and Model

1.1 Regression vs. Correlation

- 1.1.1 Roles of explanation vs. prediction in science
- 1.1.2 Response vs. explanatory variables
- 1.1.3 Regression vs. correlation: relationship betw/ variables
 - regression: one variable (explanatory) *determines* other (response)
 - correlation: both variables change together; *covary*
- 1.1.4 When to use regression or correlation?
 - regression to measure effect of one variable on other
 - correlation to measure strength of association between variables

1.2 Other kinds of regression

- 1.2.1 nonlinear regression
- 1.2.2 multiple (= multivariate) regression
- 1.2.3 logistic regression

1.3 Linear Model for Simple Linear Regression

$$Y_i = \alpha + \beta X_i + \epsilon_i \quad \begin{array}{l} X_i = \text{explanatory variable} \\ Y_i = \text{response variable} \end{array}$$

1.4 Fitting Regression Model

- 1.4.1 Least Squares estimates for a, b :
 - minimize squared deviations of data points from regression line;

$$\rightarrow \text{minimize } \sum_{i=1}^n [Y_i - (a + bX_i)]^2$$

1.4.2 $\Sigma x^2, \Sigma y^2, \Sigma xy$

$$\begin{array}{ll} \Sigma x^2 = \sum (X_i - \bar{X})^2 & \text{machine formula: } \Sigma x^2 = \sum X_i^2 - (\sum \bar{X}^2) / n \\ \Sigma y^2 = \sum (Y_i - \bar{Y})^2 & \Sigma xy = \sum (X_i - \bar{X})(Y_i - \bar{Y}) \end{array}$$

1.4.3 Regression coefficient (b)

$$b = \frac{\sum xy}{\sum x^2} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

potential values for b : $(-\infty, \infty)$

1.4.4 Y intercept (a):

from linear model, $Y_i = \alpha + \beta X_i \Rightarrow \alpha = \bar{Y} - \beta \bar{X}$

estimate: $a = \bar{Y} - b\bar{X}$

need both a, b to specify unique regression line

1.5 Assumptions

- 1.5.1 For each X value, normal distribution of Y values (\Rightarrow normal distribution of ϵ 's)
- 1.5.2 Equal variances – e.g., range, SD of Y values not \uparrow w/ larger X
- 1.5.3 Actual relationship is linear
 - (1.5.1)-(1.5.3) can be addressed w/ transformations
- 1.5.4 Y 's randomly sampled & independent of each other
- 1.5.5 X measurements w/out error (impossible; \Rightarrow assumption effectively is that X error negligible)
- 1.5.6 Concern about outliers

4 Confidence Intervals in Regression

4.1 General form for confidence intervals: $CI = \text{statistic} \pm (t_\alpha)(SE \text{ of statistic})$

4.2 Confidence interval for regression coefficient, b

$$1-\alpha \text{ confidence limits: } b \pm t_{\alpha(2),(n-2)}s_b$$

5 Interpreting Regression Results

Regression \cong fitting line to data; quantitative exercise

Does not prove that relationship exists

6 Data Transformations in Regression

6.1 Purpose of transformations: to adjust distribution of data to satisfy assumptions (i.e., normality, equality of variances)

→ not to straighten curved lines

– e.g., log or sq.root transform may straighten points into line, but then other assumptions violated

6.2 Transformations of explanatory data (X) not affect distribution of Y , so can be used to straighten curved line

Caution with transformation of response data (Y):

→ inappropriate transformation will violate assumptions

6.3 Inspection of Residuals

7 Comparing Two Slopes

7.1 Hypotheses: $H_0: \beta_1 = \beta_2$ $H_A: \beta_1 \neq \beta_2$

7.2 Student's t -test:

7.3 test statistic:

$$t = \frac{b_1 - b_2}{s_{b_1 - b_2}} \qquad s_{b_1 - b_2} = \sqrt{\frac{(s_{Y.X}^2)_p}{(\sum x^2)_1} + \frac{(s_{Y.X}^2)_p}{(\sum x^2)_2}}$$

$$(s_{Y.X}^2)_p = \frac{(\text{residual SS})_1 + (\text{residual SS})_2}{(\text{residual DF})_1 + (\text{residual DF})_2}$$

7.4 Critical value (t) has $(n_1 - 2) + (n_2 - 2)$ degrees of freedom:

$$\nu = n_1 + n_2 - 4$$

7.5 If H_0 not rejected, estimate common regression coefficient:

$$b_c = \frac{(\sum xy)_1 + (\sum xy)_2}{(\sum x^2)_1 + (\sum x^2)_2}$$

or (w/ more rounding error):

$$b_c = \frac{(\sum x^2)_1 b_1 + (\sum x^2)_2 b_2}{(\sum x^2)_1 + (\sum x^2)_2}$$