

Review on Linear Regression Analysis

Example: City X has recently conducted a travel survey of its 10 zones. The collected data are summarized in Table 1. Based on the given information,

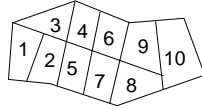


Table 1 Observed Trips and Zonal Activity

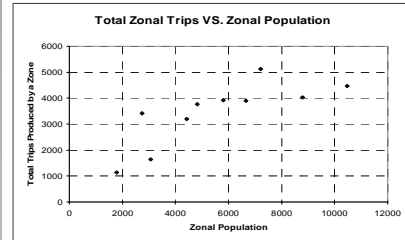
Zone #	Pop. X_1	Households X_2	Trips Y
1	7212	2488	5126
2	4818	2188	3764
3	8789	2423	4030
4	5805	2141	3921
5	3054	1241	1644
6	10463	3451	4467
7	2735	1857	3407
8	1784	905	1143
9	4418	1695	3202
10	6657	1960	3900

- How to predict the total number of trips that will be produced by each zone in 10 years (assume zonal population in 10 years is known)?
- How to relate Y (total trips produced by a zone) to X_1 (zonal population)?

Our First Attempt (informal):

City X

- To answer the questions, we need to relate Y to the X s!
- Plot the observed Y vs X_1 (scatter chart) and then Add a **Trend Line**



Outline:

- Simple Linear Regression Analysis
- Multiple Linear Regression
- Understanding the Regression Outputs
- Validating a Regression Model

Simple Linear Regression

- **Problem:** consider two variables, X and Y , and we have n paired observations on them (data): $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$, what is the relationship between Y and X ?

- **Method:** Assumed Y and X are **linearly** related, that is

$$Y = b_0 + b_1 X$$

where: X - **independent** variable (a factor that is considered to have impact on Y)

Y - **dependent** variable

b_0, b_1 - **regression coefficients**

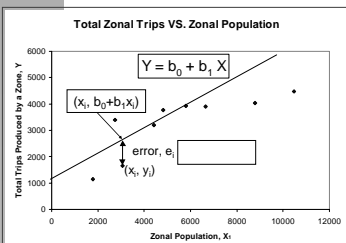
- **Our Goals:**

- Find the **best** value for b_0 and b_1 (but what do we mean by "best"?)
- Find out whether or not the identified relationship is good, i.e., validate the model

Our Second Attempt (formal):

City X

- Plot Y vs. X_1 (scatter graph) and the straight line $Y = b_0 + b_1 X_1$ with different intercept (b_0) and slope (b_1)



We shall find b_0 and b_1 so as to **minimize** the total estimation error (E):

$$E = \sum e_i^2 = \sum [y_i - (b_0 + b_1 x_i)]^2$$

This is called the **method of least squares**.

Results of Our Second Attempt:

City X

- b_0 and b_1 can be determined mathematically as follows:

Estimating the Regression Coefficients

- Data: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- From data we can calculate the following statistics

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$S_{XX} = \sum_{i=1}^n (x_i - \bar{x})^2, \quad S_{YY} = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$S_{XY} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}), \quad S_{Y|X} = \sqrt{\frac{S_{YY} - S_{XY}^2 / S_{XX}}{n-2}}$$

- b_0 and b_1 can then be determined by

$$b_0 = \bar{y} - b_1 \bar{x}, \quad b_1 = \frac{S_{XY}}{S_{XX}}$$

7

Our Third Attempt (using Excel):

SUMMARY OUTPUT

Regression Statistics

Multiple R	0.780
R ²	0.608
Adjusted R ²	0.560
Standard Error	808.122
Observations	10

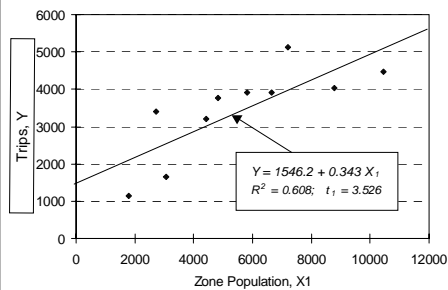
Coefficients / Standard Error / t Stat

Intercept	1546.214	600.034	2.577
X Variable 1	0.343	0.097	3.526

b_0
 b_1

8

Scatter Plot and Regression Line



9

Multiple Regression

- Example 1 (Continued): it seems that the total number of trips produced by a zone should also be related to the households of that zone (in addition to population), can we relate Y to both X_1 (zonal population) and X_2 (zonal households)?

- In general:

- Data: $(y_1, x_{11}, x_{12}, \dots), (y_2, x_{21}, x_{22}, \dots), \dots, (y_n, x_{n1}, x_{n2}, \dots)$
- Assumed Relationship

$$Y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_k x_k$$

- Problem: how to find coefficients b_0, b_1, b_2, \dots

- The *principle* to solve the problem is the same as for the simple linear regression, that is, the **method of least squares**.

10

Multiple Regression Using Excel

City X

SUMMARY OUTPUT

Regression Statistics

Multiple R	0.891
R ²	0.794
Adjusted R ²	0.735
Standard Error	627.121
Observations	10

Coefficients / Standard Error / t Stat

Intercept	491.361	637.458	0.771
X Variable 1	0.804	0.575	1.398
X Variable 2	0.552	0.398	1.385

Resulting Equation: $Y = 491.361 + 0.804 X_1 + 0.552 X_2$

11

Understanding Regression Output

- We obtained $Y = b_0 + b_1 X$, but

- how good is the final regression equation?
- How good are the estimates b_0 and b_1 ?
- ...

- Regression coefficients are estimates of the "true" population coefficients. If the true relationship is

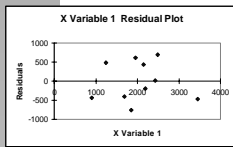
$$Y = \beta_0 + \beta_1 x$$

then b_0 and b_1 are estimates of β_0 and β_1 and a different sample would result in different b_0 and b_1 !!!

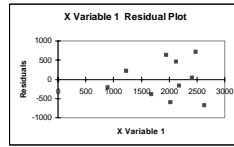
12

Check for Homogeneous Variance

- **Objective:** is the standard deviation of the residual (e) constant across all X values?
- **Method:** check scatter plot of the residuals VS. each X variable



homogeneous variance

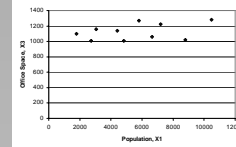


non-homogeneous variance

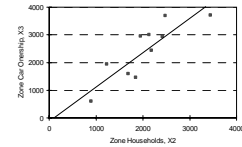
26

Check for Multicollinearity

- **Objective:** is independent variable X_i linearly correlated to independent variable X_j ?
- **Method:** check scatter plot of X_i vs. X_j



No evidence of multicollinearity



Evidence of multicollinearity

27

Check for Multicollinearity - cont.

- **Sample Correlation Coefficient (r):** is a measure of the degree of linear-relationship between two variables.
- $-1 \leq r \leq 1$:
 - if $r = +1$, X_1 and X_2 have a perfect *positive* correlation
 - if $r = -1$, X_1 and X_2 have a perfect *negative* correlation
 - if $r = 0$, X_1 and X_2 have *no* linear relationship
 - if $|r| > 0.4$ for two independent variables, it may run into the multicollinearity problem if both variables are included in a regression equation.
- Correlation Coefficient can also be directly obtained using Excel function

28