STAT 3008 Applied Regression Analysis

TUTORIAL 2: Simple Linear Regression

LAI Chun Hei*

Department of Statistics, The Chinese University of Hong Kong

1 Model Assumption

To quantify the relationship between two factors, say X and Y, we have to at least assume the type of relationship they have, could it be linear, logarithm, quadratic, etc.

In the experiment, the factor Y is likely to be affected by the factor X. In usual experiment context, X and Y are known as independent and dependent factors respectively. But in linear regression, X and Y are known as predictor and response respectively. Now, we state our model assumption:

Simple Linear Regression Model

The model relating the two factors are assumed to be:

$$y_i = Y_{|X=x_i|} = \beta_0 + \beta_1 x_i + e_i$$

where $\mathbb{E}(e_i) = 0$, $\operatorname{Var}(e_i) = \sigma^2$ and e_i 's are i.i.d. Therefore, we have

$$\mathbb{E}(Y|X = x_i) = \beta_0 + \beta_1 x$$
$$\operatorname{Var}(Y|X = x_i) = \sigma^2$$

Remark 1.1. The values x_i here are known constants, instead of some realised observations from a random variable X. For regression of random predictors, you may refer to [1].

Remark 1.2. It should be noticed that the response y_i is a random variable. The data set $\{x_i, y_i\}_{i=1}^n$ consists of <u>realised values</u> from the random variables.

2 Least Square Estimator

Essentially, we want to fit a straight line to the set of points on the Cartesian plane. However, there are many ways to define "good" in a fit. The simplest way is to consider the <u>total vertical distance</u> between the points and the line. The best fit line is therefore the line which minimises the distance, i.e. which minimises

$$RSS(\beta_0, \beta_1) = \sum_{i}^{n} [y_i - (\beta_0 + \beta_1 x_i)]^2$$

Minimising the total distance is equivalent to minimising the total squared distance. Hence, we have the Least Square Method. By elementary multivariate calculus and statistical concepts, the derivation performed in the lesson yields the following results.

^{*}s1155002282@sta.cuhk.edu.hk

Least Squares Estimator

If we define the following notations,

$$\bar{x} = \frac{1}{n} \sum_{i} x_{i}$$

$$\bar{y} = \frac{1}{n} \sum_{i} y_{i}$$

$$SXX = \sum_{i} (x_{i} - \bar{x})^{2} = \sum_{i} x_{i}^{2} - n \bar{x}^{2}$$

$$SYY = \sum_{i} (y_{i} - \bar{y})^{2} = \sum_{i} y_{i}^{2} - n \bar{y}^{2}$$

$$SXY = \sum_{i} (x_{i} - \bar{x})(y_{i} - \bar{y}) = \sum_{i} x_{i}y_{i} - n \bar{x} \bar{y}$$

then the estimators are given by

$$\begin{split} \hat{\beta_0} &= \bar{y} - \hat{\beta_1} \, \bar{x} \\ \hat{\beta_1} &= \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} = \frac{\text{SXY}}{\text{SXX}} \\ \hat{\sigma^2} &= \frac{\sum_i \hat{e_i}^2}{n - 2} = \frac{\text{SYY} - \text{SXY}^2 / \text{SXX}}{n - 2} \end{split}$$

Under our estimated model, we therefore have the fitted values given x_i and the residual, i.e. the difference between the fitted value and the realised value.

$$\hat{y}_i = \hat{\mathbb{E}}(Y|X = x_i) = \hat{\beta}_0 + \hat{\beta}_1 x_i$$
$$\hat{e}_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$$

Remark 2.1. For the detailed derivation, please refer to the lecture notes. You should be familiar with their derivation as they may be tested in midterm and final exam.

Remark 2.2. $\hat{\beta}_0$ and $\hat{\beta}_1$ can be written as linear combinations of y_i :

$$\hat{\beta}_1 = \sum_i \left(\frac{x_i - \bar{x}}{\mathrm{SXX}}\right) y_i$$
 and $\hat{\beta}_0 = \sum_i \left[\frac{1}{n} - \bar{x}\left(\frac{x_i - \bar{x}}{\mathrm{SXX}}\right)\right] y_i.$

This is useful when deriving the distribution and consistency of the estimators in Exercise 4.1.

Remark 2.3. By the derivative condition of RSS with respect to β_0 , we have

$$\sum_{i} \hat{e}_{i} = 0 \quad \text{and} \quad \bar{y} = \hat{\beta}_{0} + \hat{\beta}_{1} \,\bar{x} \,.$$

Exercise 2.1. (2012 Fall Midterm #3) Use the simple linear regression model to fit a straight line on two data points: (-2, 4), (-1, 3). What are the values of $\hat{\beta}_0$ and $\hat{\beta}_1$?

Exercise 2.2. Show that

$$\sum_{i} x_i \, \hat{e_i} = 0 \qquad \text{and, therefore} \qquad \sum_{i} \hat{y_i} \, \hat{e_i} = 0.$$

Exercise 2.3. Show that $\hat{\sigma^2}$ is an unbiased estimator of σ^2 , i.e.

 $\mathbb{E}(\hat{\sigma^2}) = \sigma^2.$

3 Analysis of Variance(ANOVA)

It is natural that we are interested in whether our model assumption is correct. The basic question will be if Y is really related to X. Mathematically, this question is equivalent to asking whether β_1 is zero in our model. The most common way is the Analysis of Variance, which compares two models of different mean functions.

We want to test the following hypothesis:

$$H_0: \mathbb{E}(Y|X=x) = \beta_0$$
 vs $H_1: \mathbb{E}(Y|X=x) = \beta_0 + \beta_1 x.$

Therefore, we need some test statistics related to these two hypothesis and, most importantly, of known distributions. Consider the residual sum of square for the two models. For H_0 , we have

$$\hat{\mathbb{E}}(Y|X=x) = \hat{\beta_0} = \bar{y} \qquad \Rightarrow \qquad \operatorname{RSS}_{H_0} = \sum_i (y_i - \bar{y})^2 = \operatorname{SYY}.$$

On the other hand, for H_1 , we have

$$\operatorname{RSS}_{H_1} = \sum_{i} \left[y_i - \left(\hat{\beta}_0 + \hat{\beta}_1 \, x_i \right) \right]^2 = \operatorname{SYY} - \frac{\operatorname{SXY}^2}{\operatorname{SXX}}.$$

Since we use more variables to fit the points in H_1 , it must be true that $\text{RSS}_{H_0} \geq \text{RSS}_{H_1}$. Therefore, H_1 is valid only when $\text{RSS}_{H_0} \gg \text{RSS}_{H_1}$. For easy comparison, we define the Sum of Square due to Regression as

$$SSreg = RSS_{H_0} - RSS_{H_1} = \frac{SXY^2}{SXX}$$

Equivalently, H_1 is valid only when SSreg $\gg 0$.

3.1 Distributions of Estimators

To define how large is large, we need the distributions as well so that we can define "large" in a probabilistic sense. Under H_0 , the distribution of the sum of squares are given below,

By simple algebra, we rewrite

$$SSreg = \frac{SXY^2}{SXX} = \left[\sum_i \left(\frac{x_i - \bar{x}}{\sqrt{SXX}}\right) y_i\right]^2$$

By Central Limit Theorem, we have

$$\sum_{i} \left(\frac{x_i - \bar{x}}{\sqrt{\text{SXX}}}\right) y_i \rightsquigarrow \mathcal{N}(0, \sigma^2).$$

Therefore, we know SSreg $/\sigma^2 \sim \chi^2(1)$. Also, it is known that $(n-2)\hat{\sigma^2}/\sigma^2 \sim \chi^2(n-2)$. The test statistic is thus given by

$$F = \frac{\frac{\text{SSreg} / \sigma^2}{1}}{\frac{(n-2)\sigma^2 H_1 / \sigma^2}{n-2}} = \frac{\text{SSreg}}{\hat{\sigma}^2 H_1} \sim F(1, n-2).$$

For significance level α , we reject H_0 if the p-value is smaller than $F_{1-\alpha}(1, n-2)$.

Remark 3.1. That $(n-2)\hat{\sigma}^2/\sigma^2 \sim \chi^2(n-2)$ is due to a fact that the degree of freedom is the number of values in the statistics that are free to vary. The detailed proof of the distribution of $(n-2)\hat{\sigma}^2/\sigma^2$ includes the use of quadratic forms, which is beyond the scope. Interested students may refer to [2].

Remark 3.2. It should be noticed that the $\hat{\sigma^2}$ here is under the model in H_1 . The reason of using this instead of $\hat{\sigma^2}$ under H_0 is examined in Exercise 3.1.

Exercise 3.1. What is the estimator of σ^2 under H_0 ? Explain why the use of it in the denominator makes no sense.

Exercise 3.2. Show that

$$\sum_{i} (y_i - \bar{y})^2 = \sum_{i} (y_i - \hat{y}_i)^2 + \sum_{i} (\hat{y}_i - \bar{y})^2$$

3.2 ANOVA Table

Thanks to the result of Exercise 3.2, we have a neat and tidy representation of ANOVA, which is called the ANOVA table.

| Source | df | \mathbf{SS} | MS | F | p-value |
|------------|-----|----------------------|---|---------------------------------------|--------------------|
| Regression | 1 | SSreg | $\mathrm{SSreg}/1$ | $\operatorname{SSreg}/\hat{\sigma^2}$ | $P(F_{1,n-1} > F)$ |
| Residual | n-1 | RSS_{H_1} | $\hat{\sigma^2} = \operatorname{RSS}/(n-2)$ | | |
| Total | n-2 | SYY | | | |

You should be extremely familiar with the above table because it appears in every midterm. We will practice this in Exercise 4.3 and 4.4.

Remark 3.3. Therefore, we can define the "Coefficient of Determination" to be

$$R^2 = \frac{\text{SSreg}}{\text{SYY}} \in [0, 1].$$

The realised value summarises the strength of relationship between the sampled response and predictors. \blacksquare

Intervals, Tests and Band 4

Besides testing the mean functions by ANOVA, we will also want to perform test on individual parameters. Therefore, we need the distributions of the estimator. We begin this section with an exercise.

Exercise 4.1. Prove that $\hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased and find their variance and asymptotic distribution.

4.1 Confidence Intervals and Tests for Intercept and Slope

With the distributions of the estimator and some facts in statistics, we can construct the test statistics from the distribution derived in Exercise 4.1.

Confidence Interval and Test for Intercept If we want to test whether the intercept is a certain value β_0^* , i.e.

$$H_0: \beta_0 = \beta_0^* \qquad \text{vs} \qquad H_1: \beta_0 \neq \beta_0^*,$$

then the test statistic is

$$t = \frac{\hat{\beta}_0 - \beta_0^*}{\operatorname{se}(\hat{\beta}_0)} \sim t(n-2) \quad \text{where} \quad \operatorname{se}(\hat{\beta}_0) = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\operatorname{SXX}}}.$$

Therefore, for significance level α , we reject H_0 when $|t| > t_{1-\frac{\alpha}{2}}(n-2)$. Also, the $(1-\alpha) \times 100\%$ confidence interval of β_0 is given by

$$\hat{\beta}_0 - t_{1-\frac{\alpha}{2}}(n-2)\operatorname{se}(\hat{\beta}_0) \le \beta_0 \le \hat{\beta}_0 + t_{1-\frac{\alpha}{2}}(n-2)\operatorname{se}(\hat{\beta}_0).$$

Confidence Interval and Test for Slope

Similarly, for the test of slope, i.e.

$$H_0: \beta_1 = \beta_1^* \qquad \text{vs} \qquad H_1: \beta_1 \neq \beta_1^*,$$

the test statistic is

$$t = \frac{\hat{\beta}_1 - \beta_1^*}{\operatorname{se}(\hat{\beta}_1)} \sim t(n-2) \quad \text{where} \quad \operatorname{se}(\hat{\beta}_0) = \frac{\hat{\sigma}}{\sqrt{\operatorname{SXX}}}.$$

Therefore, for significance level α , we reject H_0 when $|t| > t_{1-\frac{\alpha}{2}}(n-2)$. Also, the $(1-\alpha) \times 100\%$ confidence interval of β_1 is given by

$$\hat{\beta}_1 - t_{1-\frac{\alpha}{2}}(n-2)\operatorname{se}(\hat{\beta}_1) \le \beta_1 \le \hat{\beta}_1 + t_{1-\frac{\alpha}{2}}(n-2)\operatorname{se}(\hat{\beta}_1).$$

Remark 4.1. Obviously, a test of zero slope, i.e.

$$H_0: \beta_1 = 0 \qquad \text{vs} \qquad H_1: \beta_1 \neq 0,$$

is equivalent to testing

$$H_0: \mathbb{E}(Y|X=x) = \beta_0$$
 vs $H_1: \mathbb{E}(Y|X=x) = \beta_0 + \beta_1 x$

which is our ANOVA F-test in Section 3. Therefore, they should give the same result. Mathematically, if we look at the t-statistics,

$$t = \frac{\hat{\beta_1} - 0}{se(\hat{\beta_1})} = \frac{\hat{\beta_1}}{\hat{\sigma}/\sqrt{SXX}}$$
$$t^2 = \frac{\hat{\beta_1}^2}{\hat{\sigma^2}/SXX} = \frac{\hat{\beta_1}^2 SXX}{\hat{\sigma^2}} = \frac{SXY^2}{\hat{\sigma^2}SXX} = \frac{SSreg}{\hat{\sigma^2}} = F.$$

In general, we have

$$F(1,m) = \frac{\chi^2(1)}{\chi^2(m)/m} = \frac{Z^2}{\chi^2(m)/m} = \left(\frac{Z}{\sqrt{\chi^2(m)/m}}\right)^2 = t(m)^2$$

Exercise 4.2. Construct a 95% confidence interval for the slope from the data set $\{(1, 1), (4, 9), (10, 10)\}$, given $t_{0.975}(1) = 12.7062$. Bosco argues that the confidence interval you construct has a 95% probability of including the true slope. Explain whether he is correct.

Exercise 4.3. (2013 Fall Midterm #1) Fill in the missing values in the following tables of regression output from a data set of size 100.

| | ANOVA Table | | | | | | |
|-------------------|------------------------|-----------------|--------------|---------|--|--|--|
| _ | Source | df SS | MS F | | | | |
| _ | Regression | | | | | | |
| _ | Residual | | | | | | |
| _ | Total | | | | | | |
| | | | | | | | |
| Coefficient Table | | | | | | | |
| Variable | Coefficient | s.e. | t-statistics | p-value | | | |
| Constant | 0.5854 | | | 0.2188 | | | |
| X | | 0.4927 | | | | | |
| n = | $\hat{\sigma} = 4.714$ | $R^2 = 0.03294$ | | | | | |

Exercise 4.4. (2012 Spring Midterm #1) Fill in the missing values in the following tables of regression output. In R, it is found that $qf(1-9.5e^{-9}, 1, 6) = 1917.3$. Also, $\bar{x} = 5.125, \bar{y} = -9.1974, SXX = 54.875$.

| _ | ANOVA Table | | | | | | | | |
|----------------|-----------------|------------------|---------------|---------|--------------|---------|--|--|--|
| | Source | df | \mathbf{SS} | MS | F | p-value | | | |
| _ | Regression | | | | | 9.5e-09 | | | |
| _ | Residual | | | | _ | | | | |
| | Total | | | | | | | | |
| | | C | officiar | + Tabla | | | | | |
| | Coemcient Table | | | | | | | | |
| | Variable | Coefficient | s | .e. | t-statistics | p-value | | | |
| (| Constant | | | | | 0.00322 | | | |
| | Х | -2.04245 | | | | | | | |
| \overline{n} | | $\hat{\sigma} =$ | $R^2 =$ | | | | | | |

4.2 Confidence and Prediction Intervals

Besides the true intercept and slope, we are often interested in the mean of response given a predictor x_* , i.e. $\mathbb{E}(Y|X = x_*)$. This can be constructed from the fitted value \hat{y}_* since it is an unbiased estimator of the mean. For example, we consider the mean IQ of a student who scores 98 in the midterm.

On the other hand, instead of the mean, we may also be interested in the response itself, i.e. $y_{|X=x_*}$. In this case, we want to make a prediction on what the outcome will be, given x_* . Here, for example, we are looking for the IQ of a student who scores 98 in the midterm.

Therefore, we want to construct intervals for the mean and the prediction. The results are listed below.

Confidence Interval for Mean

Given a predictor value x_* , the true value and estimation of the mean are respectively

$$\mathbb{E}(Y|X = x_*) = \beta_0 + \beta_1 x_*$$
 and $\hat{y}_* = \mathbb{E}(Y|X = x_*) = \beta_0 + \beta_1 x_*.$

The estimation uncertainty of the mean is

$$\operatorname{Var}(\hat{y_*} - \mathbb{E}(Y|X = x_*)) = \operatorname{Var}(\hat{\beta_0} + \hat{\beta_1} x_*) = \sigma^2 \left(\frac{1}{n} + \frac{(x_* - \bar{x})^2}{\operatorname{SXX}}\right).$$

Define the standard error of fit as

$$\operatorname{sefit}(\hat{y_*}|x_*) = \hat{\sigma}\sqrt{\frac{1}{n} + \frac{(x_* - \bar{x})^2}{\mathrm{SXX}}} \quad \Rightarrow \quad \frac{(n-2)\hat{\sigma^2}}{\sigma^2} = \frac{(n-2)\operatorname{sefit}(\hat{y_*}|x_*)^2}{\frac{1}{n} + \frac{(x_* - \bar{x})^2}{\mathrm{SXX}}} \sim \chi^2(n-2).$$

Therefore, the $(1 - \alpha) \times 100\%$ confidence interval for the mean is given by

$$\hat{y}_* - t_{1-\frac{\alpha}{2}} \operatorname{sefit}(\hat{y}_*|x_*) \le \mathbb{E}(Y|X=x_*) \le \hat{y}_* + t_{1-\frac{\alpha}{2}} \operatorname{sefit}(\hat{y}_*|x_*).$$

Prediction Interval for Response

Given a predictor value x_* , the response and its estimation are respectively

$$y_{|X=x_*} = \beta_0 + \beta_1 x_* + e$$
 and $\hat{y}_* = \beta_0 + \beta_1 x_*.$

The estimation uncertainty of the prediction is

$$\operatorname{Var}\left(\hat{y_{*}} - y_{|X=x_{*}}\right) = \operatorname{Var}(\hat{\beta_{0}} + \hat{\beta_{1}} x_{*} + e) = \sigma^{2} \left(1 + \frac{1}{n} + \frac{(x_{*} - \bar{x})^{2}}{\operatorname{SXX}}\right)$$

Define the standard error of prediction as

$$\operatorname{sepred}(\hat{y_*}|x_*) = \hat{\sigma}\sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\mathrm{SXX}}} \quad \Rightarrow \quad \frac{(n-2)\hat{\sigma^2}}{\sigma^2} = \frac{(n-2)\operatorname{sefit}(\hat{y_*}|x_*)^2}{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\mathrm{SXX}}} \sim \chi^2(n-2).$$

Therefore, the $(1 - \alpha) \times 100\%$ confidence interval for the mean is given by

$$\hat{y}_* - t_{1-\frac{\alpha}{2}} \operatorname{sepred}(\hat{y}_*|x_*) \le y_{|X=x_*} \le \hat{y}_* + t_{1-\frac{\alpha}{2}} \operatorname{sepred}(\hat{y}_*|x_*).$$

Remark 4.2. The estimation uncertainty of prediction and mean differs only by σ^2 . The extra uncertainty comes from the error term in the new observation that we wants to predict. Compare

$$y_{|X=x} = \beta_0 + \beta_1 x + e$$
 and $\mathbb{E}(Y|X=x) = \beta_0 + \beta_1 x.$

Due to the extra uncertainty, prediction interval includes and is larger than the confidence interval for mean.

4.3 Confidence Band

In the previous subsection, we construct confidence interval of mean for a certain point x. It is tempting to connect all the upper limits and lower limits of confidence intervals, i.e.

$$(\hat{\beta}_0 + \hat{\beta}_1 x) \pm t_{1-\frac{\alpha}{2}} \operatorname{sefit}(\hat{y}|x), \quad \forall x$$

and say that this random band has a $(1 - \alpha) \times 100\%$ probability of including the true mean line $\mathbb{E}(Y|X = x) = \beta_0 + \beta_1 x$. However, this is wrong (see Remark 4.3 and Exercise 4.5). The correct band is given below.

Confidence Band for Mean Function The $(1 - \alpha) \times 100\%$ confidence band of the mean function is given by

$$\mathcal{C}(x) = (\hat{\beta}_0 + \hat{\beta}_1 x) \pm \sqrt{2F_{1-\alpha}(2, n-2)} \operatorname{sefit}(\hat{y}|x), \quad \forall x.$$

Therefore, it is true that

 $\Pr(\text{The mean line lies in the confidence band}) = \Pr\left(\forall x, \mathbb{E}(Y|X=x) \in \mathcal{C}(x)\right) = 1 - \alpha.$

Remark 4.3. For confidence interval $C(x) = (\hat{\beta}_0 + \hat{\beta}_1 x) \pm t_{1-\frac{\alpha}{2}} \operatorname{sefit}(\hat{y}|x_*)$, we have by definition

$$\forall x, \Pr\left(\mathbb{E}(Y|X=x) \in C(x)\right) = 1 - \alpha.$$

This relationship holds for each point, i.e. pointwise. While for the confidence band, we have

$$\Pr\left(\forall x, \mathbb{E}(Y|X=x) \in \mathcal{C}(x)\right) = 1 - \alpha.$$

Here, the inclusion is for the entire line. The two cases are different.

Exercise 4.5. Explain, why it is wrong to say the band,

$$(\hat{\beta}_0 + \hat{\beta}_1 x) \pm t_{1-\frac{\alpha}{2}} \operatorname{sefit}(\hat{y}|x), \quad \forall x$$

has a $(1 - \alpha) \times 100\%$ probability of including the mean line $\mathbb{E}(Y|X = x) = \beta_0 + \beta_1 x$.

.

Exercise 4.6. For the data set $\{(1, 1), (4, 9), (10, 10)\}$, construct

- 1. a 95% confidence interval and a 95% interval for the point x = 3, and
- 2. a 95% confidence band.
- 3. What is the value of the band when x = 3?

You are given $t_{0.975}(1) = 12.7062$ and $F_{0.95}(2, 1) = 199.5$.

5 Residuals

To check whether our model assumption is valid, a good way is to look at the residual plot. Recall that the residuals

$$\hat{e}_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$$

so this gives

$$\mathbb{E}(\hat{e}_i) = \mathbb{E}(y_i) - \beta_0 - \beta_1 x_i = 0 \quad \text{and} \quad \operatorname{Var}(\hat{e}_i) = \sigma^2 \left(1 + \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\operatorname{SXX}} \right)$$

The data set gives a set of <u>realised</u> \hat{e}_i . According to our observation above, these realised residuals should

- have mean close to zero, and
- have constant variance for all value x_i .

A plot that satisfies the above criteria is a null plot, which indicates that the model assumption is valid and the regression is a good fit.

6 Appendix

For more reference, you may refer to the following text books.

References

- [1] DOUGLAS C. MONTGOMERY, ELIZABETH A. PECK AND G. GEOFFREY VINING (2006). Introduction to Linear Regression Analysis, Wiley.
- [2] ROBERT V. HOGG AND ALLEN T. CRAIG Introduction to Mathematical Statistics, Pearson.