

STAT 417 Lecture Note 24

In previous lectures, we used χ^2 -test to study the relationships among discrete variables. The χ^2 -test is only useful for discrete data. When data is continuous, χ^2 -test is no longer useful. A more powerful approach is needed.

Linear Regression:

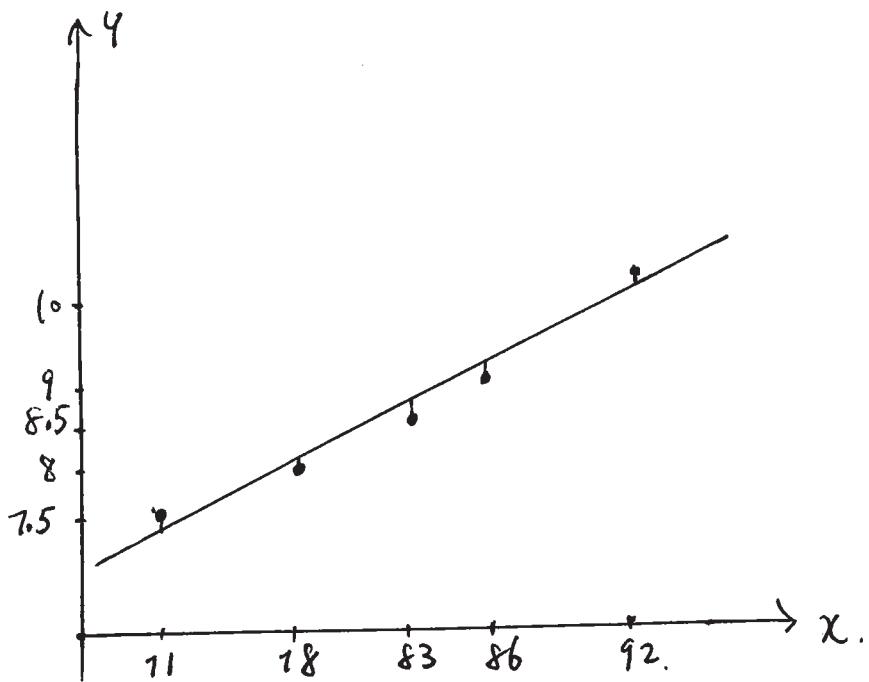
$n = 5$ students graduated last summer. Below are their grades in STAT 417 (X) and current salaries (Y) in \$10,000:

X	92	86	83	78	71
Y	10	9	8.5	8	7.5

Question: Does STAT 417 grades affect your future salary?

Both grades and salaries are continuous, so χ^2 -test is no longer useful. We will use linear regression technique to study the relationship between X and Y .

understanding the data:



from the above plot, it is easy to see that y is increasing with x . we want to establish a linear relationship between x and y . the idea is to find a line as good as possible.

$$y = b_1 + b_2 x \text{ to fit the data as good as possible.}$$

$$SE = \sum_{i=1}^5 (Y_i - b_1 - b_2 x_i)^2$$

want to find b_1, b_2 to minimize SE .

$$\frac{\partial}{\partial b_1} SE = - \sum_{i=1}^5 2(Y_i - b_1 - b_2 x_i) = 0 \Rightarrow b_1 = \frac{1}{5} \sum_{i=1}^5 (Y_i - b_2 x_i) \\ = \bar{Y} - b_2 \bar{x}$$

$$\frac{\partial}{\partial b_2} SE = - \sum_{i=1}^n 2x_i(y_i - b_1 - b_2 x_i) = 0.$$

$$\Rightarrow - \sum_{i=1}^n 2x_i(y_i - \bar{y} + b_2 \bar{x} - b_2 x_i)$$

$$= - \sum_{i=1}^n 2x_i(y_i - \bar{y} - b_2(x_i - \bar{x})) = 0$$

$$\text{so } b_2 = \frac{\sum_{i=1}^n x_i(y_i - \bar{y})}{\sum_{i=1}^n x_i(x_i - \bar{x})}$$

$$\underline{\text{Ex1: show:}} \quad \sum_{i=1}^n x_i(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

$$\sum_{i=1}^n x_i(x_i - \bar{x}) = \sum_{i=1}^n (x_i - \bar{x})^2.$$

$$\begin{aligned} \underline{\text{proof:}} \quad \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) &= \sum_{i=1}^n x_i(y_i - \bar{y}) - \sum_{i=1}^n \bar{x}(y_i - \bar{y}) \\ &= \sum_{i=1}^n x_i(y_i - \bar{y}). \end{aligned}$$

$$\text{Note: } \sum_{i=1}^n \bar{x}(y_i - \bar{y}) = \bar{x} \sum_{i=1}^n (y_i - \bar{y}) = \bar{x} \left(\sum_{i=1}^n y_i - n\bar{y} \right) = 0.$$

$$\sum_{i=1}^n \bar{x}(x_i - \bar{x}) = \bar{x} \sum_{i=1}^n (x_i - \bar{x}) = \bar{x} \left(\sum_{i=1}^n x_i - n\bar{x} \right) = 0$$

$$\text{so } \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i(x_i - \bar{x}) - \sum_{i=1}^n \bar{x}(x_i - \bar{x}) = \sum_{i=1}^n x_i(x_i - \bar{x}).$$

By Ex1, we get the following

$$b_2 = \frac{\sum_{i=1}^5 (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^5 (x_i - \bar{x})^2} = 0.118.$$

so. $b_1 = \bar{y} - b_2 \bar{x} = -1.076.$

so. the line we want is $y = -1.076 + 0.118x.$

This line is called the "line of best fit".

General case: Suppose we have n samples

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n).$$

The line of best fit, or formally, the best fitting line, is

$$y = b_1 + b_2 x.$$

where

$$\left\{ \begin{array}{l} b_2 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ b_1 = \bar{y} - b_2 \bar{x} \end{array} \right.$$

b_1, b_2 are the estimators of the intercept β_1 and slope β_2 .

prediction of future ~~salary~~ salary.

suppose the 6th graduate student has STAT 417 grade 88.
then how much she will earn?

Recall the ~~fitted~~ line of the best fit is

$$y = -1.076 + 0.118 x.$$

when $x = 88$.

$$y = 9.3$$

so ~~the~~ it is predicted that the future salary is ~~\$~~ $\$93,000$

Theorem 5.

$$\sum_{i=1}^n (y_i - \bar{y})^2 = b_2^2 \sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n (y_i - b_1 - b_2 x_i)^2.$$

Proof:

$$\begin{aligned}
\sum_{i=1}^n (y_i - b_1 - b_2 x_i)^2 &= \sum_{i=1}^n (y_i - \bar{y} - b_2(x_i - \bar{x}))^2 \\
&= \sum_{i=1}^n (y_i - \bar{y})^2 - 2 \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) b_2 + b_2^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\
&= \sum_{i=1}^n (y_i - \bar{y})^2 - 2 \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} b_2 \cdot \sum_{i=1}^n (x_i - \bar{x})^2 \\
&\quad + b_2^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\
&= \sum_{i=1}^n (y_i - \bar{y})^2 - 2 b_2 \sum_{i=1}^n (x_i - \bar{x})^2 + b_2^2 \sum_{i=1}^n (x_i - \bar{x})^2 \\
&= \sum_{i=1}^n (y_i - \bar{y})^2 - b_2^2 \sum_{i=1}^n (x_i - \bar{x})^2.
\end{aligned}$$

If y and x have a strong relationship, then b_2 should be far away from zero (i.e. b_2^2 is large). This means, $b_2^2 \sum_{i=1}^n (x_i - \bar{x})^2$ should take a large proportion of $\sum (y_i - \bar{y})^2$.

Does your future salary depend on your STAT 417 score?

An ANOVA test. ANOVA = analysis of variance.

use the following so called ANOVA Table:

Source	df	Sum of Squares	Mean squares
X	1	$b_2 \sum_{i=1}^n (x_i - \bar{x})^2$	$b_2^2 \frac{1}{2} \sum_{i=1}^n (x_i - \bar{x})^2$
Error	n-2	$\sum_{i=1}^n (y_i - b_1 - b_2 x_i)^2$	$s^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - b_1 - b_2 x_i)^2$
Total	n-1	$\sum_{i=1}^n (y_i - \bar{y})^2$	

F-test: F is named after Ron Fisher.

$$F = \frac{b_2^2 \frac{1}{2} \sum_{i=1}^n (x_i - \bar{x})^2}{s^2} \approx F(1, n-2).$$

In intuitively, a larger F value means a stronger relationship between X and Y.

The distribution of F can be used to determine if F is big or small.

In our grade-salary example, the ANOVA table is

Source	df	Sum of Squares (SS)	MS
X	1	3.54	3.54
ERROR	3	0.16	0.053
Total	4	3.7	

$$\begin{aligned} \text{so } F &= \frac{3.54}{0.053} \\ &= 65.6 \\ &> F_{0.95}(1, 3) \\ &= 10.13. \end{aligned}$$

STAT 417 Final Review:

A brief review:

1. ~~posterior~~ Bayes model: {
 - Bayesian estimation { posterior expectation
 - credible interval. (only for normal model) posterior mode
 Bayesian hypothesis testing { prior odds ratio
 posterior odds ratio
 Bayes factor. }
2. model check: {
 by χ^2 -test.
 check. normal model.
 check. multinomial model.
 check. exponential model
3. independence check: {
 by χ^2 -test
 Check independence of categorical variables among.
 Check independence among continuous variables by linear regression
4. Linear Regression: {
~~test fit line~~. $y \in b_1 + b_2 x$.
 line of the best fit: $b_1 = \dots$
 $b_2 = \dots$
~~ANOVA test prediction~~: $y_{\text{new}} = b_1 + b_2 x_{\text{new}}$.
 ANOVA Test if y and x are independent.
 i.e. test $H_0: \text{slope} = 0$