

Math 218 Mathematical Statistics

Prof. D. Joyce, Clark University

16 Mar 2009

Second Test. Wednesday, 25 Mar 2009. On chapters 6–9.

Due Today. From chapter 8, p. 290, exercises 1, 2, 3, 9.

Due Friday. From chapter 8, exercise 14, and from Chapter 9, exercises 1–3, 6.

Last time. Inferences for proportions and count data: point estimators, interval estimators, and hypothesis tests.

Today. Comparing two proportions, introduction to linear regression.

Inferences for comparing two proportions. Here we have two Bernoulli populations with unknown parameters p_1 and p_2 .

Suppose we have random samples X_1, \dots, X_{n_1} and Y_1, \dots, Y_{n_2} from these two populations. Sample means \bar{X} and \bar{Y} are point estimators \hat{p}_1 and \hat{p}_2 , respectively. These estimators have means p_1 and p_2 and variances p_1q_1/n_1 and p_2q_2/n_2 , respectively.

For large samples the difference $\hat{p}_1 - \hat{p}_2$ is approximately normally distributed with mean $p_1 - p_2$ and variance $p_1q_1/n_1 + p_2q_2/n_2$. An estimate for this variance is $\hat{p}_1\hat{q}_1/n_1 + \hat{p}_2\hat{q}_2/n_2$, so the statistic

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\hat{p}_1\hat{q}_1/n_1 + \hat{p}_2\hat{q}_2/n_2}}$$

is approximately standard normal, so it can be used to make statistical inferences about $p_1 - p_2$.

A $(1 - \alpha)$ -confidence interval for $p_1 - p_2$ has endpoints

$$\hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} \sqrt{\frac{\hat{p}_1\hat{q}_1}{n_1} + \frac{\hat{p}_2\hat{q}_2}{n_2}}.$$

As in the last chapter, we can have hypothesis tests for the difference of the means where the null hypothesis $H_0 : p_1 - p_2 = d$ is that the means differ by a constant d . Typically, however, the hypothesized difference is $d = 0$, and $H_0 : p_1 = p_2$. In that case, a pooled estimate of p , the overall probability of success, is just \hat{p} , the proportion of successes in all $n_1 + n_2$ trials, that's just a weighted sum of \hat{p}_1 and \hat{p}_2

$$\hat{p} = \frac{n_1\hat{p}_1 + n_2\hat{p}_2}{n_1 + n_2}.$$

Therefore, the statistic

$$Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}\hat{q}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

is approximately standard normal and can be used for statistical inferences.

Introduction to linear regression.

Linear regression started out as the *method of least squares*, a method developed at the beginning of the 19th century to find the closest line to a bivariate data set. The method of least squares started out with Legendre's creation of the method of least squares as published in 1806 in his work *Nouvelles méthodes pour la détermination des orbites des comètes*. (Gauss claimed he discovered earlier, but he was young at the time and had not communicated his discovery until 1809. Laplace also wrote about it in 1812.)

Here's the idea of the method of least squares. You have n data points

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

and you want to find the linear function $y = ax + b$ whose graph is closest to the points. Since we're treating y as a function of x , the error for a given point (x_i, y_i) is measured by how far off the actual value y_i is from the predicted value for the function $y = ax + b$ which would be $y = ax_i + b$. Thus, the error for that point is

$$|(ax_i + b) - y_i|.$$

Legendre's idea was to use the square of the error instead of the error itself, then add all the errors together to get the total error $\mathcal{E}(a, b)$ for the given line $y = ax + b$:

$$\mathcal{E}(a, b) = \sum_{i=1}^n (ax_i + b - y_i)^2.$$

Using standard methods from calculus, we can find the line with the smallest total error $\mathcal{E}(a, b)$ by taking the derivatives with respect to a and b and setting them to 0 to find the critical point, which, in this case there's only one (assuming $n \geq 2$ and the x values are not all the same). The minimum will occur at that critical point.

Here are the details for that computation. Note that each summation is $\sum_{i=1}^n$.

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial a} &= \sum (ax_i + b - y_i) \cdot x_i \\ &= a \sum x_i^2 + b \sum x_i - \sum x_i y_i \\ \frac{\partial \mathcal{E}}{\partial b} &= \sum (ax_i + b - y_i) \\ &= a \sum x_i + bn - \sum y_i \end{aligned}$$

When these two partial derivatives are each set to 0, we get the pair of simultaneous linear equations

$$\begin{cases} a \sum x_i^2 + b \sum x_i = \sum x_i y_i \\ a \sum x_i + bn = \sum y_i \end{cases}$$

in the two unknowns a and b .

One way to solve any pair of linear equations in two unknowns is to use Cramer's rule. It says that the pair of simultaneous linear equations

$$\begin{cases} pa + qb = r \\ sa + tb = u \end{cases}$$

$$\begin{aligned} a &= \frac{rt - uq}{pt - sq} \\ b &= \frac{pu - sr}{pt - sq} \end{aligned}$$

For our pair of equations, that gives

$$\begin{aligned} a &= \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{n \sum x_i^2 - (\sum x_i)^2} \\ b &= \frac{(\sum x_i^2)(\sum y_i) - (\sum x_i)(\sum x_i y_i)}{n \sum x_i^2 - (\sum x_i)^2} \end{aligned}$$

For statistics, it's most useful to express the answer in terms of sample means $\bar{x} = \frac{1}{n} \sum x_i$ and $\bar{y} = \frac{1}{n} \sum y_i$. Dividing the numerators and denominators of our solution equations by n gives us

$$\begin{aligned} a &= \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} \\ b &= \frac{(\sum x_i^2) \bar{y} - \bar{x} (\sum x_i y_i)}{\sum x_i^2 - n \bar{x}^2} \end{aligned}$$

If we add a little notation, namely,

$$\begin{aligned} S_{xy} &= \sum (x_i - \bar{x})(y_i - \bar{y}) \\ &= \sum x_i y_i - n \bar{x} \bar{y} \\ S_{xx} &= \sum (x_i - \bar{x})^2 \\ &= \sum x_i^2 - n \bar{x}^2 \\ S_{yy} &= \sum (y_i - \bar{y})^2 \\ &= \sum y_i^2 - n \bar{y}^2 \end{aligned}$$

then we can write our answer even more succinctly as

$$\begin{aligned} a &= \frac{S_{xy}}{S_{xx}} \\ b &= \bar{y} - a \bar{x} \end{aligned}$$

Thus, we have found the *least squares line*. It is the line

$$y = ax + b$$

where a and b are as just found.