

- Last time :
- Sample size  $\leq 10$ , can use sample range
  - usually, use  $s^2$  to estimate  $\sigma^2$
  - chi-square distribution
  - confidence interval for  $\sigma^2$ :  $\left[ \frac{(n-1)s^2}{\chi^2_{\alpha/2}}, \frac{(n-1)s^2}{\chi^2_{1-\alpha/2}} \right]$
  - hypothesis testing  $\chi^2 = \frac{(n-1)s^2}{\sigma^2}$

If we have a two tailed test, we can use a confidence interval to perform the test.

- e.g. Statistics Canada claims that the standard deviation in the heights of Canadian women is 5cm. We randomly select 11 women, and put a sample standard deviation at 6cm. Test the claim at a .05 level of significance.

$$H_0: \sigma^2 = 5^2$$

$$H_1: \sigma^2 \neq 5^2$$

The  $100(1-\alpha)\% = 95\%$  confidence interval

$$\text{For } \sigma^2 \text{ is } \left[ \frac{10(6)^2}{\chi^2_{0.025}}, \frac{10(6)^2}{\chi^2_{0.975}} \right]$$

$$\Rightarrow \left[ \frac{10(6)^2}{20.983}, \frac{10(6)^2}{3.247} \right] = (17.58, 110.87)$$

As  $s^2$  lies in the 95% confidence interval, we cannot reject the claim at a .05 level of confidence.

- Suppose we have two normal populations, and we want to compare their variances,  $\text{Pop}_1 = \sigma_1^2$   
 $\text{Pop}_2 = \sigma_2^2$

If we make the assumption that the variances are the ... and we take a random sample of size  $n_1$  from  $\text{Pop}_1$ ,  $n_2$  from  $\text{Pop}_2$ , and let  $F = \frac{s_1^2}{s_2^2}$

Then  $F$  has distribution with  $\nu_1 = n_1 - 1$ ,  $\nu_2 = n_2 - 1$

We have  $F_{.01}$  and  $F_{.05}$  values.

$$\Pr(F > F_{\alpha}) = \alpha$$

$$\left. \begin{array}{l} H_0: \sigma_1^2 \leq \sigma_2^2 \\ H_1: \sigma_1^2 > \sigma_2^2 \end{array} \right\} F = \frac{s_1^2}{s_2^2}, \text{ reject that if } F > F_{\alpha}$$

$$\left. \begin{array}{l} H_0: \sigma_1^2 \geq \sigma_2^2 \\ H_1: \sigma_1^2 < \sigma_2^2 \end{array} \right\} \text{ swap pops 1 + 2}$$

$H_0: \sigma_1^2 = \sigma_2^2$   
 $H_1: \sigma_1^2 \neq \sigma_2^2$

IF necessary, swap so that  $S_1^2 \geq S_2^2$   
 Repeat that if  $F > F_{\alpha/2}$

- e.g. Krusty claims that the standard deviation in the number of marshmallows in boxes of Lucky charms is at least as great as that in Krusty-O's. We randomly select 10 boxes of Lucky charms and 8 boxes of Krusty O's

For the lucky charms, we get a sample standard deviation of 10, for the Krusty O's, 15.

Test the claim at a .05 level of significance.

Pop<sub>1</sub> = Krusty-O's

Pop<sub>2</sub> = Lucky Charms

$H_0: \sigma_1^2 \leq \sigma_2^2$

$$F = \frac{S_1^2}{S_2^2}$$

and reject that if  $F > F_{.05}$

$H_1: \sigma_1^2 > \sigma_2^2$

$$F = \frac{15^2}{10^2} = 2.25$$

As  $z_1 = 7, z_2 = 9, F_{.05} = 3.29$

As  $F \not> F_{.05}$ , we cannot reject the claim at a .05 level of significance.

## Chapter 11 - Regression Analysis

Let  $x$  and  $y$  be random variables

$y$  will depend upon  $x$  plus a randomness factor

$y = f(x) + E$ , when  $f$  is a function and  $E$  is a random variable with  $\text{non-0}$ . We call  $f(x)$  the regression curve. Finding it is called regression analysis.

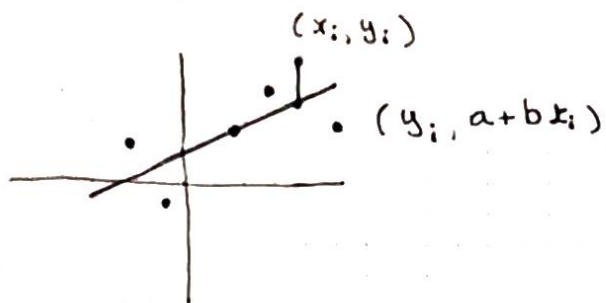
We can get an idea by taking a random sample  $(x_1, y_1), \dots, (x_n, y_n)$ , and drawing a scatterplot.

In linear regression, we have  $y = \alpha + \beta x + E$ ,  $\alpha, \beta \in \mathbb{R}$

We will find the line of best fit for our data,

$$\hat{y} = a + bx$$

We will use the method of least squares



We want to minimize the sum of the squares of the vertical distance from the points to the line

$$\sum_{i=1}^n (y_i - (a + bx_i))^2$$

$$a : \sum_{i=1}^n 2(y_i - (a + bx_i))(-1) = 0$$

$$\sum_{i=1}^n a + \sum_{i=1}^n bx_i = \sum_{i=1}^n y_i$$

$$na + (\sum_{i=1}^n x_i)b = \sum_{i=1}^n y_i \quad (1)$$

$$b : \sum_{i=1}^n 2(y_i - (a + bx_i))(-x_i) = 0$$

$$\sum_{i=1}^n ax_i + \sum_{i=1}^n bx_i^2 = \sum_{i=1}^n x_i y_i$$

$$(\sum_{i=1}^n x_i)a + (\sum_{i=1}^n x_i^2)b = \sum_{i=1}^n x_i y_i \quad (2)$$

Solve (1), (2) For a, b

- e.g. Find the line of best fit for:

$x_i$	1	2	3
$y_i$	7	3	1

$$\rightarrow na + \sum x_i b = \sum y_i$$

$$(1) \quad 3a + 6b = 11$$

$$\rightarrow \sum x_i a + \sum x_i^2 b = \sum x_i y_i$$

$$(2) \quad 6a + 14b = 16$$

$$(2) - 2(1) \Rightarrow 2b = -6, \quad b = -3$$

$$\text{Subbing into (1): } 3a + 6(-3) = 11, \quad a = \frac{29}{3}$$

$$\hat{y} = \frac{29}{3} - 3x$$

- in the above example, estimating when

$$x = 2.5$$

$$\hat{y} = \frac{29}{3} - 3(2.5)$$

- Find the line of best fit for:

$x_i$	1	2	3	4
$y_i$	2	4	5	7

$$\rightarrow na + \sum x_i b = \sum y_i$$

$$\textcircled{1} \quad 4a + 10b = 18$$

$$\rightarrow \sum x_i a + \sum x_i^2 b = \sum x_i y_i$$

$$\textcircled{2} \quad 10a + 30b = 53$$

$$\textcircled{2} - 3\textcircled{1} : \quad -2a = -1$$

$$a = 1/2$$

$$\text{Subbing into } \textcircled{1} : \quad 4(1/2) + 10b = 18$$

$$b = 1.6$$

$$\hat{y} = (0.5) + (1.6)x$$



- e.g. we have a normal population with  $\sigma = 10$   
Suppose we wish to test.

$$H_0: \mu \leq 100$$

$$H_1: \mu > 100$$

we will take a random sample of size 25 and reject  $H_0$  if  $\bar{x} > 104$   
Find  $\alpha$ .

We assume  $\mu = 100$

$$\begin{aligned} \Pr(\bar{x} > 104) &= 1 - \Pr(\bar{x} \leq 104) \\ &= 1 - F\left(\frac{104 - 100}{10/\sqrt{25}}\right) = 1 - F(2) = 1 - 0.9772 \\ &= 0.0228 \end{aligned}$$

- e.g. binomial: bernoulli trials } replacement!  
(Flipping coin, rolling die)

- e.g. At least 4 hearts, and 1 seven?

4 hearts, including 7  
 $\binom{12}{3} 36$

4 hearts, not including 7  
 $\binom{12}{4} 3$

5 hearts  
 $\binom{12}{4}$

$$\Rightarrow \frac{\binom{12}{3} 36 + \binom{12}{4} 3 + \binom{12}{4}}{\binom{52}{5}}$$

- e.g. Two kings and one club?

King of clubs, another king  
 $3 \binom{36}{4}$

no king of clubs  
 $\binom{3}{2} \binom{12}{1} \binom{36}{2}$

$$\Rightarrow \frac{3 \binom{36}{4} + \binom{3}{2} 12 \binom{36}{2}}{\binom{52}{5}}$$

- e.g. A bag contains 40 red marbles and 60 blue marbles. We reach into the bag and pull out 10 marbles. Find the prob that we get 3 red marbles.

$$\frac{\binom{80}{3} \binom{60}{7}}{\binom{100}{10}}$$

$$\left\{ \begin{array}{l} N = 100 \\ n = 10 \\ G = 40 \\ x = 3 \end{array} \right.$$

w/ replacement:  $\binom{10}{3} \left(\frac{80}{100}\right)^3 \left(\frac{60}{100}\right)^7$

- e.g. Roll a balanced die 20 times, count the threes

(i) Find the prob. of at most 4 threes

(ii) Find the prob. of at least 4 threes

(i)  $B(4; 20, 1/6)$

(ii)  $1 - B(3, 20, 1/6)$

- e.g. Krusty claims that the average box of Krusty-O's contains at least 60 marshmallows. We randomly select 10 boxes and put a sample mean of 56 and a sample standard deviation of 10.

Test the claim at a .05 level of significance

$H_0: \mu \geq 60$  As  $\sigma$  is unknown,  $n < 30$

$H_1: \mu < 60$   $t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$ , reject  $H_0$  if  $t < t_{.05}$

$t = \frac{56 - 60}{10/\sqrt{10}} = -1.6$ , As  $n = 10$   
 $t_{.05} = 1.753$

As  $t \not< -t_{.05}$  we cannot reject the claim.

- e.g. Krusty claims that the standard deviation in the number of jagged metal Krusty O's per box is 5. We randomly select 41 boxes and get a standard deviation of 3.

Test the claim at a .05 level of significance

$H_0: \sigma^2 = 5^2$   $\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$ , reject if  $\chi^2 > \chi^2_{.025}$

$H_1: \sigma^2 \neq 5^2$

$\chi^2 = \frac{40(3)^2}{5^2} = 14.4$

As  $n = 41$ ,  $\chi^2_{.025} = 59.342$

$\chi^2_{.975} = 24.433$

As  $\chi^2 = 14.4 < \chi^2_{.975}$ , we reject  $H_0$

- e.g. Find the line of best fit:

$x_i$	1	2	3	4
$y_i$	6	5	3	1

$$na + \sum x_i b = \sum y_i$$

$$4a + 10b = 15 \quad (1)$$

$$\sum x_i a + \sum x_i^2 b = \sum x_i y_i$$

$$10a + (1+4+9+16)b = (6+10+9+4)$$

$$10a + 30b = 29 \quad (2)$$

Solving (1) with (2)

$$(2) - 3(1) = -2a = -16$$

$$a = 8$$

$$\text{then, } b = -1.4$$

$$\hat{y} = 8 - 1.4x$$

- Last time - Confidence intervals for two-handed test

- F-test for two variables

- regression analysis  $y = \beta(x) + E$

-  $y = \alpha + \beta x + E$

- line of best fit:  $\hat{y} = a + bx$

$$na + \sum x_i b = \sum y_i$$

$$\sum x_i a + \sum x_i^2 b = \sum x_i y_i$$

If our data is  $(x_1, y_1), \dots, (x_n, y_n)$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2, S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2, S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$b = S_{xy} / S_{xx} \Rightarrow a = \bar{y} - b\bar{x}$$

- e.g.

$x_i$	1	2	3
$y_i$	8	3	1

$$\bar{x} = \frac{1+2+3}{3} = 2, S_{xx} = ((1-2)^2 + (2-2)^2 + (3-2)^2) = 2$$

$$\bar{y} = \frac{8+3+1}{3} = 4, S_{yy} = ((8-4)^2 + (3-4)^2 + (1-4)^2) = 26$$

$$S_{xy} = (1-2)(8-4) + (2-2)(3-4) + (3-2)(1-4) = -4 + 0 - 3 = -7$$

$$b = S_{xy} / S_{xx} \Rightarrow (-7) / (2) = -7/2; a = 4 - (-7/2)(2) \Rightarrow 4 + 7 = 11$$

$$\hat{y} = 11 - \frac{7}{2}x$$

Exponential regression:  $y = \alpha \beta^x + E$

Our best fit curve will be  $\hat{y} = ab^x$

$$\log \hat{y} = \log a + x \log b$$

Let  $c = \log a$ ,  $d = \log b$ , we have  $\log \hat{y} = c + dx$

We perform linear regression with  $(x_i, \log y_i)$

$x_i$	1	2	3
$y_i$	1	10	1000
$\log y_i$	0	1	3

$$nc + \sum x_i d = \sum \log y_i$$

$$3c + bd = 4 \quad (1)$$

$$\sum x_i c + \sum x_i^2 d = \sum x_i \log y_i$$

$$6c + 14d = 11 \quad (2)$$

$$(2) - 2(1): 2d = 3, d = 3/2$$

$$\text{then } c = -5/3$$

$$a = 10^c = 10^{-5/3}, b = 10^d = 10^{3/2}$$

$$\hat{y} = 10^{-5/3} (10^{3/2})^x$$



Power regression:  $y = \alpha x^{\beta} + \epsilon$

Our best fit curve:  $\hat{y} = ax^b$

$$\log \hat{y} = \log C + b \log x. \text{ Let } C = \log$$

$$\log \hat{y} = c + d \log x \quad \text{Perform linear regression with } (\log x_i, \log y_i)$$

$x_i$	10	100	1000
$y_i$	1000	100	1
$\log x_i$	1	2	3
$\log y_i$	3	2	0

$$nc + \sum \log x_i d = \sum \log y_i$$

$$3c + 6d = 5 \quad (1)$$

$$\sum \log x_i c + \sum \log x_i^2 d = \sum (\log x_i)(\log y_i)$$

$$6c + 14d = 7 \quad (2)$$

$$(2) - 2(1): 2d = -3; d = -3/2, 3c + 6(-3/2) = 5, c = 14/3$$

$$a = 10^c \Rightarrow 10^{14/3}; b = d = -3/2$$

$$\hat{y} = 10^{14/3} x^{-3/2}$$

Reciprocal regression:  $y = \frac{1}{\alpha + \beta x} + \epsilon$

Our best fit curve is  $\hat{y} = \frac{1}{a + bx}$

$$\frac{1}{\hat{y}} = a + bx. \text{ Perform linear regression with } (x_i, 1/y_i)$$

- e.g.

$x_i$	1	2	3	4
$y_i$	1	1/3	1/4	1/5
$1/y_i$	1	3	4	5

$$na + \sum x_i b = \sum 1/y_i$$

$$4a + 10b = 13 \quad (1)$$

$$\sum x_i a + \sum x_i^2 b = \sum \frac{x_i}{y_i}$$

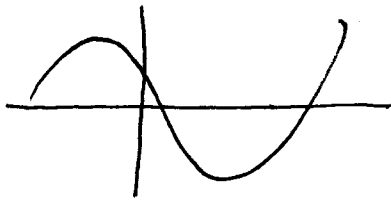
$$10a + 30b = 39 \quad (2)$$

$$3(1) - (2): 2a = 0, a = 0; b = 39/30 = 13/10$$

$$\hat{y} = \frac{1}{(13/10)x}$$

Polynomial regression:  $y = \beta_0 + \beta_1 x + \beta_2 x^2, \dots, \beta_p x^p + \epsilon$

Our best fit curve is  $\hat{y} = S_0 + S_1 x + S_2 x^2 + \dots + S_p x^p$



We will minimize the sum of the sequence of the vertical distance from the points to the curve.

$$\sum_{i=1}^n (y_i - (b_0 + b_1 x_i + b_2 x_i^2 + \dots + b_p x_i^p))^2$$

Differentiate  $b_0$  :  $\sum_{i=1}^n 2(y_i - (b_0 + b_1 x_i + b_2 x_i^2 + \dots + b_p x_i^p))(-1) = 0$

$$nb_0 + \left(\sum_{i=1}^n x_i\right)b_1 + \left(\sum_{i=1}^n x_i^2\right)b_2 + \dots + \left(\sum_{i=1}^n x_i^p\right)b_p = \sum_{i=1}^n y_i \quad (1)$$

Differentiate  $b_1$  :  $\sum_{i=1}^n 2(y_i - (b_0 + b_1 x_i + b_2 x_i^2 + \dots + b_p x_i^p))(-x_i) = 0$

$$\left(\sum_{i=1}^n x_i\right)b_0 + \left(\sum_{i=1}^n x_i^2\right)b_1 + \left(\sum_{i=1}^n x_i^3\right)b_2 + \dots + \left(\sum_{i=1}^n x_i^{p+1}\right)b_p = \sum_{i=1}^n x_i y_i \quad (2)$$

$$\left(\sum_{i=1}^n x_i^2\right)b_0 + \left(\sum_{i=1}^n x_i^3\right)b_1 + \left(\sum_{i=1}^n x_i^4\right)b_2 + \dots + \left(\sum_{i=1}^n x_i^{p+2}\right)b_p = \sum_{i=1}^n x_i^2 y_i \quad (3)$$

⋮

$$\left(\sum_{i=1}^n x_i^p\right)b_0 + \left(\sum_{i=1}^n x_i^{p+1}\right)b_1 + \left(\sum_{i=1}^n x_i^{p+2}\right)b_2 + \dots + \left(\sum_{i=1}^n x_i^{2p}\right)b_p = \sum_{i=1}^n x_i^p y_i \quad (p+1)$$

Solve for  $b_0, b_1, \dots, b_p$

-e.g. Find the quadratic of best fit for :

$x_i$	1	2	3	4
$y_i$	-2	0	3	10

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$$

$$\hat{y} = b_0 + b_1 x + b_2 x^2$$

$$nb_0 + \sum x_i b_1 + \sum x_i^2 b_2 = \sum y_i$$

$$4b_0 + 10b_1 + 30b_2 = 11 \quad (1)$$

$$\sum x_i b_0 + \sum x_i^2 b_1 + \sum x_i^3 b_2 = \sum x_i y_i$$

$$10b_0 + 30b_1 + 100b_2 = 47 \quad (2)$$

$$\sum x_i^2 b_0 + \sum x_i^3 b_1 + \sum x_i^4 b_2 = \sum x_i^2 y_i$$

$$30b_0 + 100b_1 + 354b_2 = 185 \quad (3)$$