

# Covariance

find  $\delta q$  if  $\delta x$  and  $\delta y$  are not independent

$N$  pairs of data  $(x_1, y_1), \dots, (x_N, y_N)$

$x_1, \dots, x_N \rightarrow \bar{x}$  and  $\sigma_x$

$y_1, \dots, y_N \rightarrow \bar{y}$  and  $\sigma_y$

$$\begin{aligned}
 & q_i = g(x_i, y_i) \\
 & q_1, \dots, q_N \rightarrow \bar{q} \text{ and } \sigma_q \\
 & \rightarrow q_i \approx g(\bar{x}, \bar{y}) + \frac{\partial g}{\partial x} (x_i - \bar{x}) + \frac{\partial g}{\partial y} (y_i - \bar{y})
 \end{aligned}$$

$$\begin{aligned}
 \bar{q} &= \frac{1}{N} \sum_{i=1}^N q_i \\
 &= \frac{1}{N} \sum_{i=1}^N \left[ g(\bar{x}, \bar{y}) + \frac{\partial g}{\partial x} (x_i - \bar{x}) + \frac{\partial g}{\partial y} (y_i - \bar{y}) \right]
 \end{aligned}$$

$$\sum (x_i - \bar{x}) = 0 \Rightarrow \underline{\bar{q} = g(\bar{x}, \bar{y})}$$

$$\begin{aligned}
 \sigma_q^2 &= \frac{1}{N} \sum (q_i - \bar{q})^2 \\
 &= \frac{1}{N} \sum \left[ \frac{\partial g}{\partial x} (x_i - \bar{x}) + \frac{\partial g}{\partial y} (y_i - \bar{y}) \right]^2 \\
 &= \left( \frac{\partial g}{\partial x} \right)^2 \frac{1}{N} \sum (x_i - \bar{x})^2 + \left( \frac{\partial g}{\partial y} \right)^2 \frac{1}{N} \sum (y_i - \bar{y})^2 \\
 &\quad + 2 \frac{\partial g}{\partial x} \frac{\partial g}{\partial y} \frac{1}{N} \sum (x_i - \bar{x})(y_i - \bar{y})
 \end{aligned}$$

$\sigma_q$  for arbitrary  $\sigma_x$  and  $\sigma_y$

$\sigma_x$  and  $\sigma_y$  can be correlated  $\longrightarrow$

covariance  $\sigma_{xy}$   $\longrightarrow$

$$\sigma_q^2 = \left( \frac{\partial g}{\partial x} \right)^2 \sigma_x^2 + \left( \frac{\partial g}{\partial y} \right)^2 \sigma_y^2 + 2 \frac{\partial g}{\partial x} \frac{\partial g}{\partial y} \sigma_{xy}$$

$$\sigma_{xy} = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

when  $\sigma_x$  and  $\sigma_y$  are independent  $\sigma_{xy} = 0 \longrightarrow$

$$\sigma_q^2 = \left( \frac{\partial g}{\partial x} \right)^2 \sigma_x^2 + \left( \frac{\partial g}{\partial y} \right)^2 \sigma_y^2$$

# Coefficient of Linear Correlation

$N$  pairs of values  $(x_1, y_1), \dots, (x_N, y_N)$

$$y = A + Bx$$



do  $N$  pairs of  $(x_i, y_i)$  satisfy a linear relation ?

$$r = \frac{\sigma_{xy}}{\sigma_x \cdot \sigma_y}$$
$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$



linear correlation coefficient  
or correlation coefficient

$$-1 \leq r \leq 1$$

Suppose  $(x_i, y_i)$  all lie exactly  
on the line  $y = A + Bx$

$$y_i = A + Bx_i$$

$$\bar{y} = A + B\bar{x}$$

$$y_i - \bar{y} = B(x_i - \bar{x})$$

$$r = \frac{B \sum (x_i - \bar{x})^2}{\sqrt{\sum (x_i - \bar{x})^2 \cdot B^2 \sum (x_i - \bar{x})^2}} = \frac{B}{|B|} = \pm 1$$

Suppose, there is no relationship  
between  $x$  and  $y$

$$\sum (x_i - \bar{x})(y_i - \bar{y}) \rightarrow 0$$

$$r = 0$$

if  $r$  is close to  $\pm 1$

when  $x$  and  $y$  are linearly correlated

if  $r$  is close to 0

when there is no relationship between  $x$  and  $y$   
 $x$  and  $y$  are uncorrelated

# Quantitative Significance of $r$

|                |    |    |    |     |    |    |    |    |     |    |
|----------------|----|----|----|-----|----|----|----|----|-----|----|
| Student $i$    | 1  | 2  | 3  | 4   | 5  | 6  | 7  | 8  | 9   | 10 |
| Homework $x_i$ | 90 | 60 | 45 | 100 | 15 | 23 | 52 | 30 | 71  | 88 |
| Exam $y_i$     | 90 | 71 | 65 | 100 | 45 | 60 | 75 | 85 | 100 | 80 |

calculate correlation coefficient

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

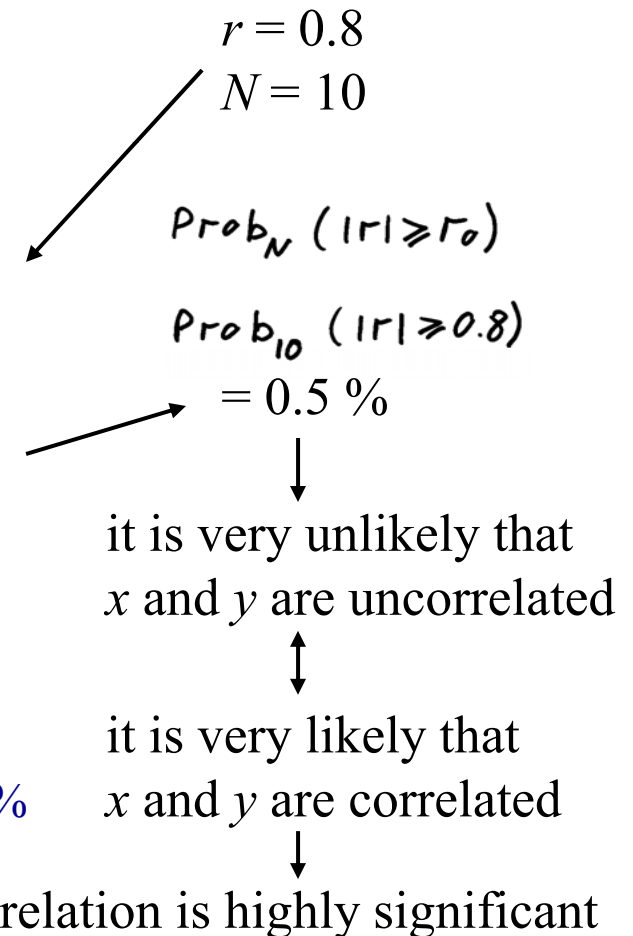
probability that  $N$  measurements of two uncorrelated variables  $x$  and  $y$  would produce  $r \geq r_0$   $\longrightarrow$  **Table C**

**Table 9.4.** The probability  $Prob_N(|r| \geq r_0)$  that  $N$  measurements of two uncorrelated variables  $x$  and  $y$  would produce a correlation coefficient with  $|r| \geq r_0$ . Values given are percentage probabilities, and blanks indicate values less than 0.05%.

| $N$ | $r_0$ |     |     |     |     |     |     |     |     |     |   |
|-----|-------|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
|     | 0     | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1 |
| 3   | 100   | 94  | 87  | 81  | 74  | 67  | 59  | 51  | 41  | 29  | 0 |
| 6   | 100   | 85  | 70  | 56  | 43  | 31  | 21  | 12  | 6   | 1   | 0 |
| 10  | 100   | 78  | 58  | 40  | 25  | 14  | 7   | 2   | 0.5 | 0   | 0 |
| 20  | 100   | 67  | 40  | 20  | 8   | 2   | 0.5 | 0.1 |     | 0   | 0 |
| 50  | 100   | 49  | 16  | 3   | 0.4 |     |     |     |     | 0   | 0 |

correlation is “significant” if  $Prob_N(|r| \geq r_0)$  is less than 5 %

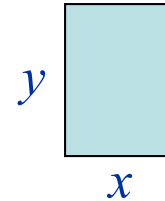
correlation is “highly significant” if  $Prob_N(|r| \geq r_0)$  is less than 1 %



Example:

Calculate the covariance and the correlation coefficient  $r$  for the following six pairs of measurements of two sides  $x$  and  $y$  of a rectangle. Would you say these data show a significant linear correlation coefficient? Highly significant?

|       | A  | B  | C  | D  | E  | F  |    |
|-------|----|----|----|----|----|----|----|
| $x =$ | 71 | 72 | 73 | 75 | 76 | 77 | mm |
| $y =$ | 95 | 96 | 96 | 98 | 98 | 99 | mm |



$$\bar{x} = 74 \quad \bar{y} = 97$$

covariance  $\sigma_{xy} = \frac{1}{N} \sum (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{6} ((-3) \times (-2) + \dots + 3 \times 2) = \underline{3}$

correlation coefficient  $r = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} = \underline{0.98}$

**Table C**  $Prob_6(|r| \geq 0.98) \approx 0.2\%$

therefore, the correlation is both significant and highly significant

| $N$ | $r_o$ |     |     |     |     |     |     |     |     |     |           |
|-----|-------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----------|
|     | 0     | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1         |
| 3   | 100   | 94  | 87  | 81  | 74  | 67  | 59  | 51  | 41  | 29  | 0         |
| 6   | 100   | 85  | 70  | 56  | 43  | 31  | 21  | 12  | 6   | 1   | <b>x0</b> |
| 10  | 100   | 78  | 58  | 40  | 25  | 14  | 7   | 2   | 0.5 |     | 0         |
| 20  | 100   | 67  | 40  | 20  | 8   | 2   | 0.5 | 0.1 |     |     | 0         |
| 50  | 100   | 49  | 16  | 3   | 0.4 |     |     |     |     |     | 0         |



# Chi Squared Test for a Distribution

40 measured values of  $x$  (in cm)

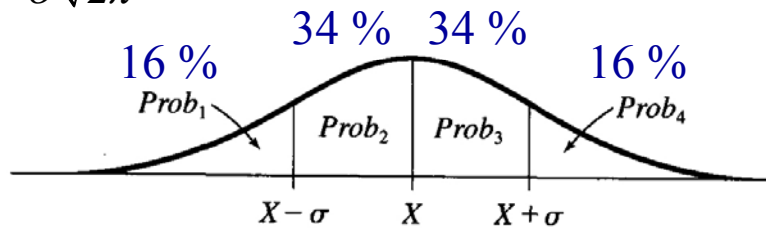
|     |     |     |     |     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 731 | 772 | 771 | 681 | 722 | 688 | 653 | 757 | 733 | 742 |
| 739 | 780 | 709 | 676 | 760 | 748 | 672 | 687 | 766 | 645 |
| 678 | 748 | 689 | 810 | 805 | 778 | 764 | 753 | 709 | 675 |
| 698 | 770 | 754 | 830 | 725 | 710 | 738 | 638 | 787 | 712 |

are these measurements governed by a Gauss distribution ?

$$\bar{x} = \frac{\sum x_i}{N} = 730.1 \text{ cm}$$

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N-1}} = 46.8 \text{ cm}$$

$$G_{X,\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-X)^2/2\sigma^2}$$



$$\frac{O_k - E_k}{\sqrt{E_k}} = \frac{\text{deviation}}{\text{expected size of fluctuation}} \sim 1 ?$$

$$\chi^2 = \sum_{k=1}^n \frac{(O_k - E_k)^2}{E_k}$$

chi squared

$\chi^2 \lesssim n$  observed and expected distributions agree about as well as expected

$\chi^2 \gg n$  significant disagreement between observed and expected distributions

| Bin number $k$ | Observed number $O_k$ | Expected number $E_k = N \text{Prob}_k$ | Difference $O_k - E_k$ |
|----------------|-----------------------|---|------------------------|
| 1              | 8                     | 6.4                                     | 1.6                    |
| 2              | 10                    | 13.6                                    | -3.6                   |
| 3              | 16                    | 13.6                                    | 2.4                    |
| 4              | 6                     | 6.4                                     | -0.4                   |

$O_k$  – observed number

$E_k$  – expected number

$\sqrt{E_k}$  – fluctuations of  $E_k$

$$\begin{aligned} \chi^2 &= \sum_{k=1}^4 \frac{(O_k - E_k)^2}{E_k} \\ &= \frac{(1.6)^2}{6.4} + \frac{(-3.6)^2}{13.6} + \frac{(2.4)^2}{13.6} + \frac{(-0.4)^2}{6.4} \\ &= 1.80 < n \end{aligned}$$

no reason to doubt that the measurements were governed by a Gauss distribution

# Degrees of Freedom and Reduced Chi Squared

a better procedure is to compare  $\chi^2$  not with the number of bins  $n$  but instead with the number of degree of freedom  $d$

$n$  is the number of bins

$c$  is the number of parameters that had to be calculated from the data to compute the expected numbers  $E_k$

$c$  is called the number of constrains

$d$  is the number of degrees of freedom

$$\underline{d = n - c}$$

test for a GAUSS distribution  $G_{\mu, \sigma}(x) \rightarrow c = 3$   $\begin{matrix} \swarrow \mu \\ \leftarrow \sigma \\ \swarrow \sigma \end{matrix}$

(expected average value of  $\chi^2$ ) =  $d = n - c$

$\tilde{\chi}^2 = \chi^2 / d$  reduced chi squared

(expected average value of  $\tilde{\chi}^2$ ) = 1

# Probabilities of Chi Squared

quantitative measure of agreement between observed data and their expected distribution

$$(\text{expected average value of } \chi^2) = d = n - c$$

$$\tilde{\chi}^2 = \chi^2 / d$$

$$(\text{expected average value of } \tilde{\chi}^2) = 1$$

$$\chi^2 = 1.80$$

$$d = 4 - 3 = 1$$

$$\tilde{\chi}^2 = 1.80$$

$$\text{Prob}(\tilde{\chi}^2 \geq 1.80) \approx 18\% \quad \leftarrow \text{Table D}$$

| d  | $\tilde{\chi}_0^2$ |      |     |      |     |      |     |      |    |     |     |     |     |
|----|--------------------|------|-----|------|-----|------|-----|------|----|-----|-----|-----|-----|
|    | 0                  | 0.25 | 0.5 | 0.75 | 1.0 | 1.25 | 1.5 | 1.75 | 2  | 3   | 4   | 5   | 6   |
| 1  | 100                | 62   | 48  | 39   | 32  | 26   | 22  | 19   | 16 | 8   | 5   | 3   | 1   |
| 2  | 100                | 78   | 61  | 47   | 37  | 29   | 22  | 17   | 14 | 5   | 2   | 0.7 | 0.2 |
| 3  | 100                | 86   | 68  | 52   | 39  | 29   | 21  | 15   | 11 | 3   | 0.7 | 0.2 | —   |
| 5  | 100                | 94   | 78  | 59   | 42  | 28   | 19  | 12   | 8  | 1   | 0.1 | —   | —   |
| 10 | 100                | 99   | 89  | 68   | 44  | 25   | 13  | 6    | 3  | 0.1 | —   | —   | —   |
| 15 | 100                | 100  | 94  | 73   | 45  | 23   | 10  | 4    | 1  | —   | —   | —   | —   |

probability of obtaining a value of  $\tilde{\chi}^2$  greater or equal to  $\tilde{\chi}_0^2$ , assuming the measurements are governed by the expected distribution

disagreement is “significant” if  $\text{Prob}_N(\tilde{\chi}^2 \geq \tilde{\chi}_0^2)$  is less than 5 %

disagreement is “highly significant” if  $\text{Prob}_N(\tilde{\chi}^2 \geq \tilde{\chi}_0^2)$  is less than 1 %

reject the expected distribution