

## Lecture 7

### Correlation example

$i$	1	2	3	4	5	6
$x$	1	2	3	4	10	10
$y$	1	3	3	5	1	11

Formula for calculating correlation

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

Alternative Formula

$$r = \frac{1}{n-1} \frac{1}{s_x} \frac{1}{s_y} \left( \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} \right)$$

Where

$r$  - correlation

$n$  - number of observations

$\bar{x}$  - sample mean of  $x$  values

$\bar{y}$  - sample mean of  $y$  values

$s_x$  - sample standard deviation of  $x$  values

$s_y$  - sample standard deviation of  $y$  values

Calculate the components of correlation formula

$$n = 6$$

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n} = \frac{1+2+3+4+10+10}{6} = \frac{30}{6} = 5$$

$$\bar{Y} = \frac{\sum_{i=1}^n Y_i}{n} = \frac{1+3+3+5+1+11}{6} = \frac{24}{6} = 4$$

sample means

Next we need the standard deviations

$$S_x = \sqrt{\frac{\sum_{i=1}^n X_i^2 - n\bar{X}^2}{n-1}}$$

← formula ~~for~~  
for standard deviation

$$S_y = \sqrt{\frac{\sum_{i=1}^n Y_i^2 - n\bar{Y}^2}{n-1}}$$

Already have  $\bar{Y}, \bar{X}, n$   
from above. Need  $\sum X^2, \sum Y^2$

$$\begin{aligned} \sum_{i=1}^n X_i^2 &= 1^2 + 2^2 + 3^2 + 4^2 + 10^2 + 10^2 \\ &= 1 + 4 + 9 + 16 + 100 + 100 \\ &= 230 \end{aligned}$$

$$\begin{aligned} \sum_{i=1}^n Y_i^2 &= 1^2 + 3^2 + 3^2 + 5^2 + 1^2 + 11^2 \\ &= 1 + 9 + 9 + 25 + 1 + 121 \\ &= 166 \end{aligned}$$

So

$$s_x = \sqrt{\frac{230 - 6(5)^2}{6-1}} = \sqrt{\frac{230 - 150}{5}} = \sqrt{16} = 4$$

$$s_y = \sqrt{\frac{166 - 6(4)^2}{6-1}} = \sqrt{\frac{166 - 96}{5}} = \sqrt{14} = 3.7416 \text{ (4dp)}$$

Lastly we need  $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$  (or  $\sum x_i y_i - n\bar{x}\bar{y}$ )

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) &= (1-5)(1-4) \\ &+ (2-5)(3-4) \\ &+ (3-5)(3-4) \\ &+ (4-5)(5-4) \\ &+ (10-5)(1-4) \\ &+ (10-5)(11-4) \\ &= (-4)(-3) + (-3)(-1) \\ &+ (-2)(-1) + (-1)(1) \\ &+ (5)(-3) + (5)(7) \\ &= 12 + 3 + 2 - 1 - 15 + 35 \\ &= \underline{36} \end{aligned}$$

$$\begin{aligned} \sum x_i y_i &= 1 \times 1 + 2 \times 3 + 3 \times 3 \\ &+ 4 \times 5 + 10 \times 1 + 10 \times 11 \\ &= 1 + 6 + 9 \\ &+ 20 + 10 + 110 \\ &= \underline{156} \end{aligned}$$

$$\begin{aligned} \sum x_i y_i - n\bar{x}\bar{y} &= 156 - 6(5)(4) \\ &= 156 - 120 \\ &= \underline{36} \end{aligned}$$

↑  
this method  
is quicker for  
hand calculation

Check that these  
agree

Now substitute back into correlation formula

$$r = \frac{1}{n-1} * \frac{1}{s_x} * \frac{1}{s_y} * \left( \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} \right)$$

$$= \frac{1}{6-1} * \frac{1}{4} * \frac{1}{3.7416} * 36$$

$$= 0.48 \quad (2dp)$$

We would say that x and y are weakly ~~and~~ positively correlated.

---

## Probability

Probability is based on the idea of a random experiment. That is an experiment

for which we cannot predict the outcome before it is carried out.



A phenomenon is random if the outcome cannot be ~~is~~ predicted but it has a regular distribution in very many repetitions of the experiment.



The probability of an event is the proportion of times that it occurs in very many repetitions of a random experiment.

In class experiment <sup>random</sup>  
we will carry out ~~a~~ an experiment using thumb tacks

Experimental procedure:

1. Drop a thumb tack from about 6 inches above a flat surface
2. Record whether it lands UP   
or down 
3. Repeat many times to estimate the probability of the thumb tack landing up.

# How to record results

Use a table. Suppose we get the following sequence of tosses: u, u, D, u, D, u, ...

then  $\Rightarrow$

Toss #	Cumulative # up	Current estimate of $P(\text{UP})$
1	1	$1/1 = 1$
2	2	$2/2 = 1$
3	2	$2/3 = 0.66$
4	3	$3/4 = 0.75$
5	3	$3/5 = 0.6$
6	4	$4/6 = 0.66$
⋮	⋮	⋮
⋮	⋮	⋮
⋮	⋮	⋮
⋮	⋮	⋮
⋮	⋮	⋮
n	$n_{up}$	$\frac{n_{up}}{n} = P(\text{UP})$

many trials {