

Last time

Sec 6.4 Covariance and the Variance of Sums

Let  $X, Y$  be RVs

$$\begin{aligned}\text{Cov}(X, Y) &= E((X - \mu_X)(Y - \mu_Y)) \\ &= E(XY) - E(X)E(Y)\end{aligned}$$

Covariance relates to correlation between  $X, Y$  as we will see today.

If  $X, Y$  are independent,  $\text{Cov}(X, Y) = 0$ , and

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

Covariance has bilinearity properties:

$$\text{Cov}(X, ZY) = Z\text{Cov}(X, Y)$$

$$\text{Cov}(X, Y+Z) = \text{Cov}(X, Y) + \text{Cov}(X, Z)$$

If  $X_1, \dots, X_n$  are exchangeable (i.e. i.d.)

$$\text{Var}(X_1 + \dots + X_n) = n\text{Var}(X_1) + n(n-1)\text{Cov}(X_1, X_2)$$

If we draw  $n$  tickets,  $X_1, X_2, \dots, X_n$  from a box with  $N$  tickets, with variance  $\sigma^2$ ,

$$\text{Cov}(X_1, X_2) = -\frac{\sigma^2}{N-1}$$

$$\text{and } \text{Var}(X_1 + \dots + X_n) = n\sigma^2 \left[ \frac{N-n}{N-1} \right]$$

↳ correction factor

Today

① review concept test from last time.

② Sec 6.4 Correlation

## Stat 134

Wednesday April 24 2019

1. Consider a Poisson( $\lambda$ ) process. Let  $T_r \sim \text{gamma}(r, \lambda)$  be the rth arrival time.  $\text{Cov}(T_1, T_3)$  equals:

- a  $\lambda$
- b  $\lambda^2$
- c  $1/\lambda^2$
- d none of the above

Recall  $\text{Var}(T_r) = \frac{r}{\lambda^2}$

Discuss with your neighbor for 1 minute how you did this.

d

Each arrival in a poisson process is independent, therefore  $T_1$  and  $T_3$  are independent which means the covariant must be zero.

c

Break  $T_3$  into  $T_3 = T_1 + W_2 + W_3$ . Then  $W_2, W_3$  are independent of  $T_1$ . By bilinearity of covariance and covariance of independent variables is 0, we can simplify to  $\text{Var}(T_1)$  so we get  $1/\lambda^2$

c

Find the  $\text{Cov}(T_1, T_3 - T_1)$ , use Cov properties and  $\text{Var}(T_1)$  to find  $\text{Cov}(T_1, T_3)$

$$\begin{aligned} 0 &= \text{Cov}(T_1, T_3 - T_1) = \text{Cov}(T_1, T_3) - \text{Var}(T_1) \\ \Rightarrow \text{Cov}(T_1, T_3) &= \lambda^2 \end{aligned}$$

## Sec 6.4 Correlation

$$\text{Cov}(x, y) = E((x - \mu_x)(y - \mu_y))$$

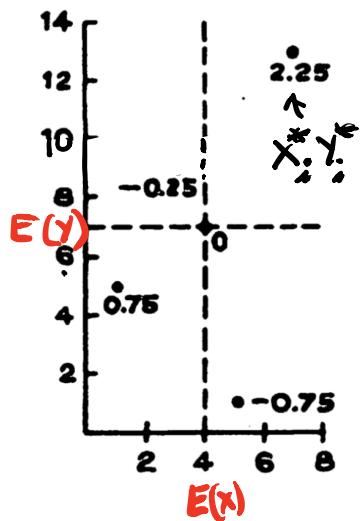
$$r = \text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\text{SD}(x)\text{SD}(y)} = E\left(\left(\frac{x - \mu_x}{\text{SD}_x}\right)\left(\frac{y - \mu_y}{\text{SD}_y}\right)\right)$$

$$= E(x^* y^*)$$

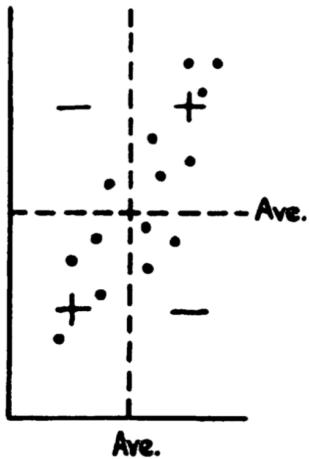
$\nwarrow$   $x, y$  in standard units

How the correlation coefficient works.

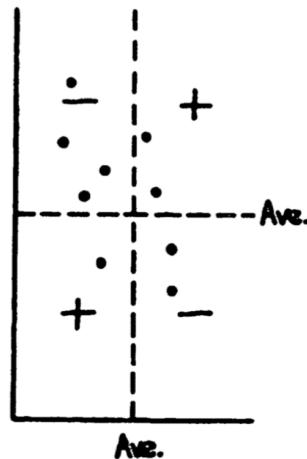
(a) Scatter diagram



(b) Positive r



(c) Negative r



This will have a positive correlation since more of the points are in the 1<sup>st</sup> and 3<sup>rd</sup> quadrants

Ex Consider a Poisson process.

$T_r \sim \text{Gamma}(r, \lambda)$  be the  $r^{\text{th}}$  arrival time.

We showed  $\text{Cov}(T_1, T_3) = \frac{1}{\lambda^2}$ .  $\left( \text{recall } \text{Var}(T_r) = \frac{r}{\lambda^2} \right)$

$$\text{Corr}(T_1, T_3) = \frac{\text{Cov}(T_1, T_3)}{\text{SD}(T_1) \text{SD}(T_3)} = \frac{\frac{1}{\lambda^2}}{\frac{1}{\lambda} \cdot \frac{\sqrt{3}}{\lambda}} = \frac{1}{\sqrt{3}}$$

$\uparrow$   
makes sense positively correlated,

Ex Suppose the sum of  $k$  exchangeable (i.e. identically distributed) RVs is a constant

$$N_1 + N_2 + \dots + N_k = c$$

Find  $\text{Corr}(N_1, N_2)$ .

sln

$$N_1 + N_2 + \dots + N_k = c$$

$$\Rightarrow \text{Var}(N_1 + \dots + N_k) = 0$$

$$\Rightarrow k\text{Var}(N_1) + k(k-1)\text{Cov}(N_1, N_2) = 0$$

$$\Rightarrow \text{Cov}(N_1, N_2) = -\frac{\text{Var}(N_1)}{k-1} \quad \frac{-\text{Var}(N_1)}{k-1}$$

$$\Rightarrow \text{Corr}(N_1, N_2) = \frac{\text{Cov}(N_1, N_2)}{\sqrt{\text{SD}(N_1)\text{SD}(N_2)}} = \boxed{-\frac{1}{k-1}}$$

//  $\text{Var}(N_1)$

## Stat 134

Friday April 26 2019

1. An urn contains 90 marbles, of which there are 20 green, 25 black, and 45 red marbles. Alice randomly picks 10 marbles without replacement, and Bob randomly picks another 10 marbles without replacement. Let  $X_1$  be the number of green marbles that Alice has and  $X_2$  the number of green marbles that Bob has.

To find  $\text{Corr}(X_1, X_2)$  is

$$X_1 + X_2 + \cdots + X_9 = 20$$

a true identity that is useful? Explain.

a yes

b no

c not enough info to decide

$$X_i = \# \text{green at } i^{\text{th}} \text{ person}.$$

$$X_i \sim \text{Hyper}(n=10, G=20, N=90) \text{ for } i=1, 2, \dots, 9.$$

Imagine randomly laying out 90 marbles flat in a line



If you pull out 10 marbles without looking at other marbles  $P(X_1 = k) = \binom{20}{k} \binom{70}{10-k}$

$$P(X_2 = k) = \frac{\binom{20}{k} \binom{70}{10-k}}{\binom{90}{10}}$$

Sum of exchangeable  $X_i$  is a constant.

2. An urn contains 90 marbles, of which there are 20 green, 25 black, and 45 red marbles. Alice randomly picks 10 marbles without replacement, and Bob randomly picks another 10 marbles without replacement. Let  $X_1$  be the number of green marbles that Alice has and  $X_2$  the number of green marbles that Bob has.

Find  $\text{Corr}(X_1, X_2)$ .

$$\begin{aligned} & \text{a } -1/8 & \Rightarrow \text{Corr}(X_1, X_2) = \frac{-1}{q-1} = \frac{-1}{8} \\ & \text{b } -1/9 \\ & \text{c } -1/10 \\ & \text{d } \text{none of the above} \end{aligned}$$

If asked to find  $\text{Cov}(X_1, X_2)$

$$\text{Cov}(X_1, X_2) = \text{Cov}(X_1, X_2) \text{SD}(X_1) \text{SD}(X_2) = \text{Cov}(X_1, X_2) \text{Var}(X_1)$$

$$X_1 \sim \text{Hyper}(10, 20, 90)$$

$$\begin{aligned} \Rightarrow \text{Var}(X_1) &= n \frac{6}{N} \frac{N-n}{N} \left[ \frac{N-n}{N-1} \right] \\ &= 10 \left( \frac{20}{90} \right) \left( \frac{70}{90} \right) \left( \frac{90-10}{90-1} \right) \end{aligned}$$

$$\Rightarrow \text{Cov}(X_1, X_2) = -\frac{1}{8} \cdot 10 \left( \frac{20}{90} \right) \left( \frac{70}{90} \right) \left( \frac{90-10}{90-1} \right) = \frac{-1400}{7209}$$

## Properties of correlation

① Correlation is invariant to change of scale except possibly by a sign.

(i.e.  $|\text{corr}(x, y)| = |\text{corr}(ax+b, cy+d)|$   
for constants  $a, b, c, d$ .

e.g. Correlation between Boston and NYC temperatures is the same whether temps in  $^{\circ}\text{C}$  or  $^{\circ}\text{F} = 1.8^{\circ}\text{C} + 32$

Proof

$$\begin{aligned}\text{corr}(ax+b, cy+d) &= \frac{\text{cov}(ax+b, cy+d)}{\text{SD}(ax+b)\text{SD}(cy+d)} \\ &= ac \frac{\text{cov}(x, y)}{\sqrt{|a||c|} \text{SD}(x)\text{SD}(y)} \\ &= \frac{ac \text{cov}(x, y)}{|a||c| \text{SD}(x)\text{SD}(y)} = \frac{ac}{|a||c|} \text{corr}(x, y)\end{aligned}$$

□

Hence

$$\text{corr}(x, y) = \text{corr}(x^*, y^*) \text{ since}$$

$$\text{SD}(x) > 0 \text{ and } \text{SD}(y) > 0.$$

$$\textcircled{2} \quad -1 \leq \text{corr}(x, y) \leq 1$$

Proof

Correlation is invariant if you convert  $x, y$  to standard units  $x^*, y^*$  since  $SD(x) > 0$ ,  $SD(y) > 0$ .

So we show that  $-1 \leq \text{corr}(x^*, y^*) \leq 1$ .

$$\left. \begin{array}{l} E(x^*) = 0 = E(y^*) \\ SD(x^*) = 1 = SD(y^*) \\ E(x^{*2}) = 1 = E(y^{*2}) \end{array} \right\}$$

Since  $x^*, y^*$   
are standard  
units.

$$E(x^{*2}) = \text{var}(x^*) + (E(x^*))^2$$

$$(x^* + y^*)^2 \geq 0$$

$$\text{so } E((x^* + y^*)^2) \geq 0$$

$$E(x^{*2} + y^{*2} + 2x^*y^*) \geq 0$$

$$1 + 1 + 2E(x^*y^*) \geq 0$$

$$E(x^*y^*) \geq -1$$

$$\Rightarrow -1 \leq E(x^*y^*)$$

$$\boxed{-1 \leq \text{corr}(x, y)}$$

Show  $\text{Corr}(x, y) \leq 1$  by examining  $E((x^* - y^*)^2)$ :

$$(x^* - y^*)^2 \geq 0$$

$$\text{so } E((x^* - y^*)^2) \geq 0$$

$$E(x^{*2} + y^{*2} - 2x^*y^*) \geq 0$$

$$1 + 1 - 2E(x^*y^*) \geq 0$$

$$E(x^*y^*) \leq 1$$

$$\Rightarrow \boxed{\text{Corr}(x, y) \leq 1}$$

This finishes the proof.

□