

Lecture 22

(1)

Continuing our discussion on jointly distributed random variables

Expected Value

Let X, Y be jointly distributed r.v. and $h(x, y)$ be an arbitrary function of X, Y .

Then

$$E[h(X, Y)] = \sum_x \sum_y h(x, y) p(x, y) \quad (X, Y \text{ are discrete})$$

$$E[h(X, Y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} h(x, y) f(x, y) dy dx \quad (X, Y \text{ are continuous})$$

Covariance is a measure of how strongly two variables are related to each other.

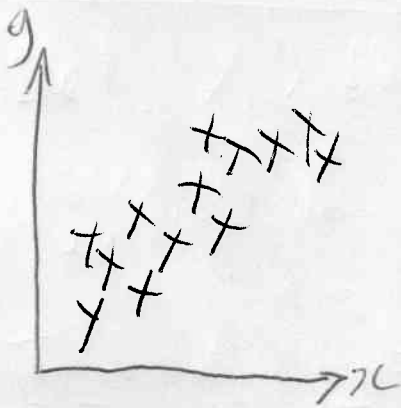
$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

12

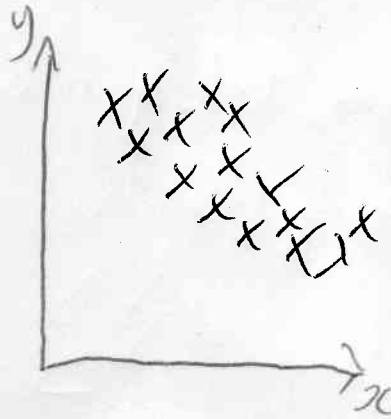
②

$$\text{Cov}(X, Y) = \sum_x \sum_y (x - \mu_x)(y - \mu_y) p(x, y) \quad X, Y \text{ discrete}$$

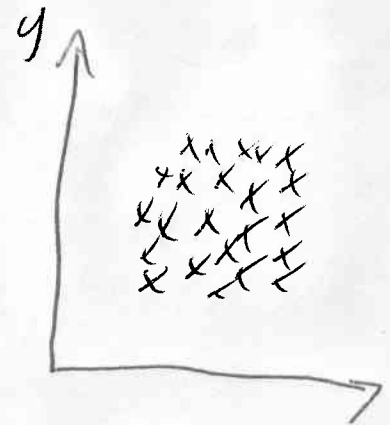
$$\text{Cov}(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu_x)(y - \mu_y) f(x, y) dy dx \quad X, Y \text{ continuous}$$



+ve
Covariance



-ve
Covariance



Covariance
near 0

Note that for computational purposes it is easier to use

$$\text{Cov}(X, Y) = E(XY) - \mu_x \mu_y$$

Proof (for continuous case)

$$\begin{aligned} \text{Cov}(X, Y) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - \mu_x)(y - \mu_y) f(x, y) dy dx \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (xy - \mu_x y - \mu_y x + \mu_x \mu_y) f(x, y) dy dx \end{aligned}$$

$$\begin{aligned}
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f(x,y) dy dx - \mu_x \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f(x,y) dy dx \\
 &\quad - \mu_y \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f(x,y) dy dx + \mu_x \mu_y \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x,y) dy dx
 \end{aligned}$$

$$= E[XY] - 2\mu_x\mu_y + \mu_x\mu_y$$

$$= E[XY] - \mu_x\mu_y$$

Unfortunately covariance depends on the units of measurement. Therefore changing the units of measurements can have considerable effect on the magnitude of the covariance. To

remedy this problem we use a rescaled version of the covariance.

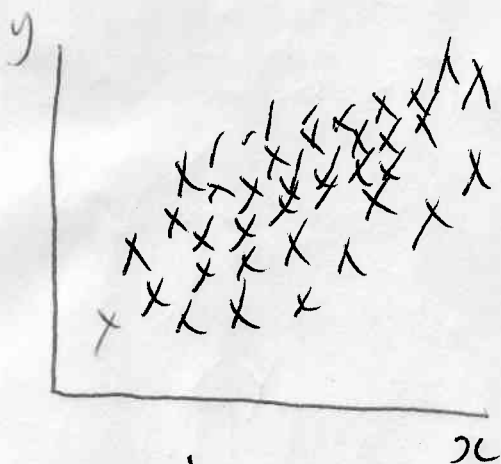
Correlation

$$\rho_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y} \quad (\text{ie we rescale by dividing by the two standard deviations})$$

Note that the following are properties of correlation

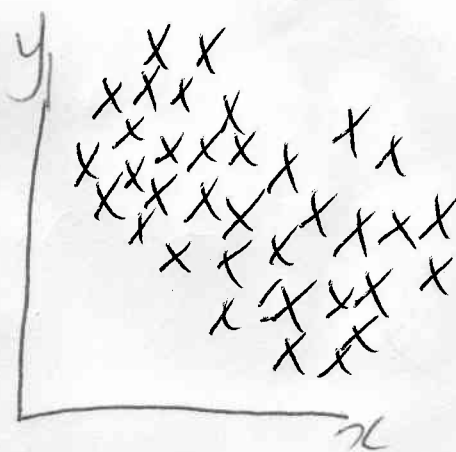
- 1. $-1 \leq \rho_{X,Y} \leq 1$ for any r.v X, Y
- 2. $\rho_{aX+b, cY+d} = \rho_{X,Y}$ (ie rescaling has no effect on correlation)
- 3. If X, Y are independent $\rho_{X,Y} = 0$
- 4. $\rho = 1$ or $-1 \iff Y = aX + b$ for some a, b ($a \neq 0$)

Interpreting correlation



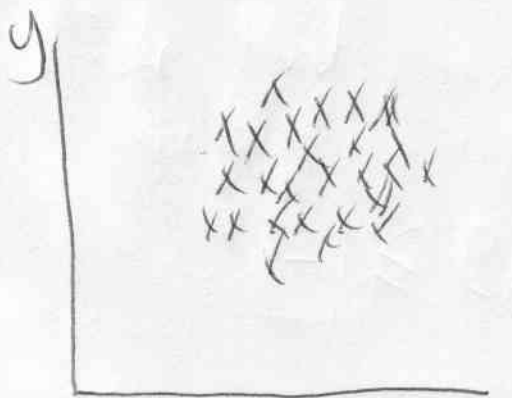
positive relationship

$$\rho_{X,Y} > 0$$



negative relationship

$$\rho_{X,Y} < 0$$



No linear relationship
 $\rho_{x,y} \approx 0$

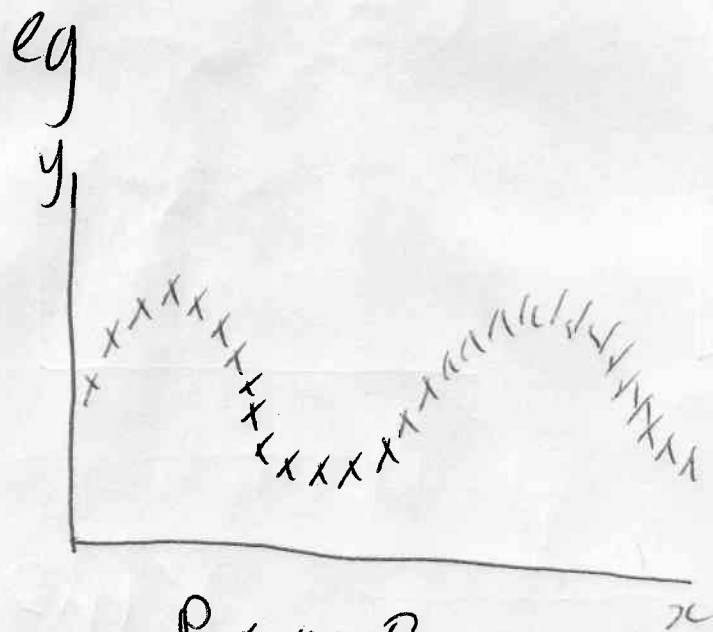


$\rho_{x,y}$ near +1
 Strong positive
 linear relationship



$\rho_{x,y}$ near -1
 Strong negative
 linear relationship

Note that correlation is only a measure of the strength of the linear relationship between x, y . It is possible to have a very strong non-linear relationship between x, y and have $\rho_{x,y} = 0$



$\rho_{x,y} = 0$

Examples

(6)

1. Recall from last time. Joint distribution of X, Y

$$f(x, y) = k(x^2 + y^2) \quad \begin{array}{l} 20 \leq x \leq 30 \\ 20 \leq y \leq 30 \end{array}$$

$$= 0 \quad \text{o.w.} \quad k = \frac{3}{380000}$$

What is $\text{Cov}(X, Y)$?

Recall that $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$

What is $E[X]$?

$$E[X] = \int_{20}^{30} x f_X(x) dx$$

$$= k \int_{20}^{30} x \left[10x^2 + \frac{30^3 - 20^3}{3} \right] dx$$

$$= k \left[\frac{10x^4}{4} + \frac{30^3 - 20^3}{6} x^2 \right]_{20}^{30}$$

$$= k \left[\frac{10(30^4 - 20^4)}{4} + \frac{(30^3 - 20^3)(30^2 - 20^2)}{6} \right]$$

$$= 25.3289$$

Similarly it can be shown that

$$E[Y] = 25.3289$$

$$E[XY] = k \int_{20}^{30} \int_{20}^{30} xy [x^2 + y^2] dx dy$$

$$= k \int_{20}^{30} \left[\frac{x^3 y^2}{2} + \frac{x y^4}{4} \right]_{20}^{30} dx$$

$$= k \int_{20}^{30} \left[\frac{x^3 (30^2 - 20^2)}{2} + x \frac{(30^4 - 20^4)}{4} \right] dx$$

$$= k \left[\frac{x^4 (30^2 - 20^2)}{8} + \frac{x^2 (30^4 - 20^4)}{8} \right]_{20}^{30}$$

$$= k \left[\frac{(30^4 - 20^4)(30^2 - 20^2)}{8} + \frac{(30^2 - 20^2)(30^4 - 20^4)}{8} \right]$$

$$= \frac{3}{380000} [81250000] = 641.4474$$

So $\text{COV}(X, Y) = 641.4474 - (25.3289)(25.3289)$
 $= -1078 \quad (4dp)$

What is $\text{Cor}(X, Y)$?

(8)

$$\text{Cor}(X, Y) = \rho_{X, Y} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

$$\text{Var}(X) = E(X^2) - (E(X))^2 \quad \text{Var}(Y) = E(Y^2) - (E(Y))^2$$

$$E(X^2) = \int_{20}^{30} x^2 f_X(x) dx$$

$$= K \int_{20}^{30} x^2 \left[10x^2 + \frac{30^3 - 20^3}{3} \right] dx$$

$$= K \left[\frac{10x^5}{5} + \frac{30^3 - 20^3}{9} x^3 \right]_{20}^{30}$$

$$= \frac{3}{380000} \left[10(30^5 - 20^5) + \frac{(30^3 - 20^3)^2}{9} \right]$$

$$= 649.8246$$

And similarly it can be shown

$$E[Y^2] = 649.8246$$

So

$$\begin{aligned}\text{Var}(X) &= 649.8246 - (25.3289)^2 \\ &= 8.2694\end{aligned}$$

$$\text{and } \text{Var}(Y) = 8.2694$$

$$\begin{aligned}\Rightarrow \sigma_x = \text{SD}(X) &= \sqrt{\text{Var}(X)} = \sqrt{8.2694} \\ &= 2.8757\end{aligned}$$

$$\sigma_y = 2.8757$$

$$\begin{aligned}\text{So } \text{Cor}(X, Y) &= \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y} = \frac{-1.078}{(2.8757)(2.8757)} \\ &= -0.130 \quad (4 \text{dp})\end{aligned}$$

So there is almost no linear relationship between x and y . (note that does not mean a strong non linear relationship might exist).