

LECTURE 12: Sums of independent random variables; Covariance and correlation

- The PMF/PDF of $X + Y$ (X and Y independent)
 - the discrete case
 - the continuous case
 - the mechanics
 - the sum of independent normals
- Covariance and correlation
 - definitions
 - mathematical properties
 - interpretation

The distribution of $X + Y$: the discrete case

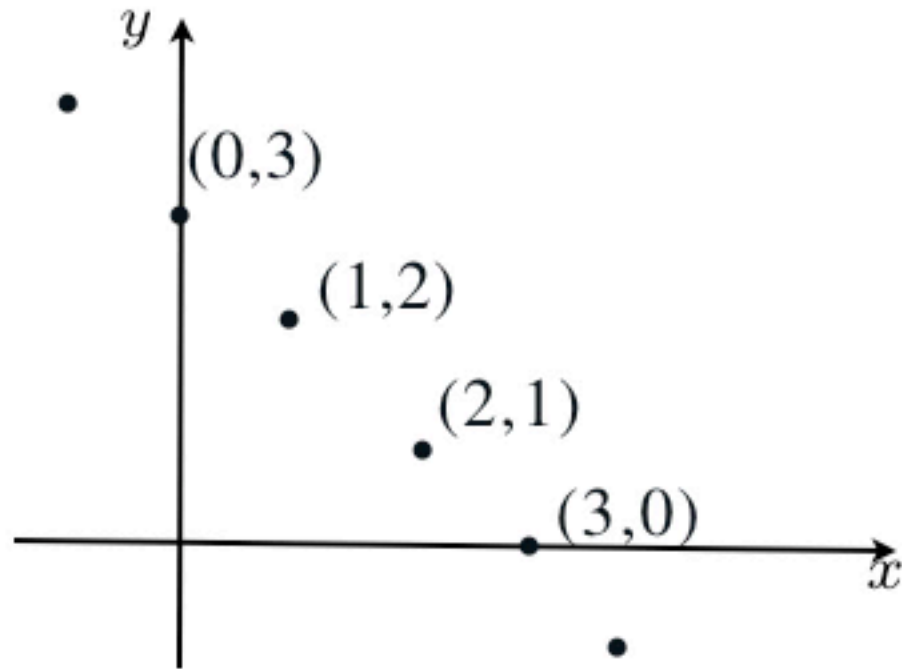
- $Z = X + Y$; X, Y independent, discrete

$g(x, y)$

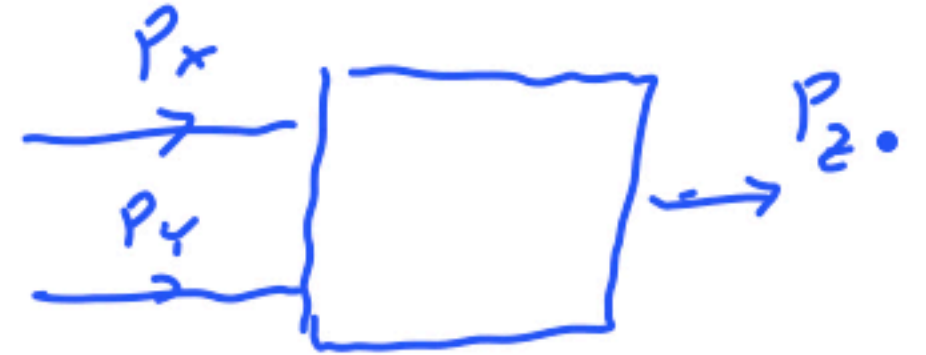
known PMFs

$$p_Z(3) = \dots + P(X=0, Y=3) + P(X=1, Y=2) + \dots$$

$$= \dots + p_X(0) p_Y(3) + p_X(1) p_Y(2) + \dots$$



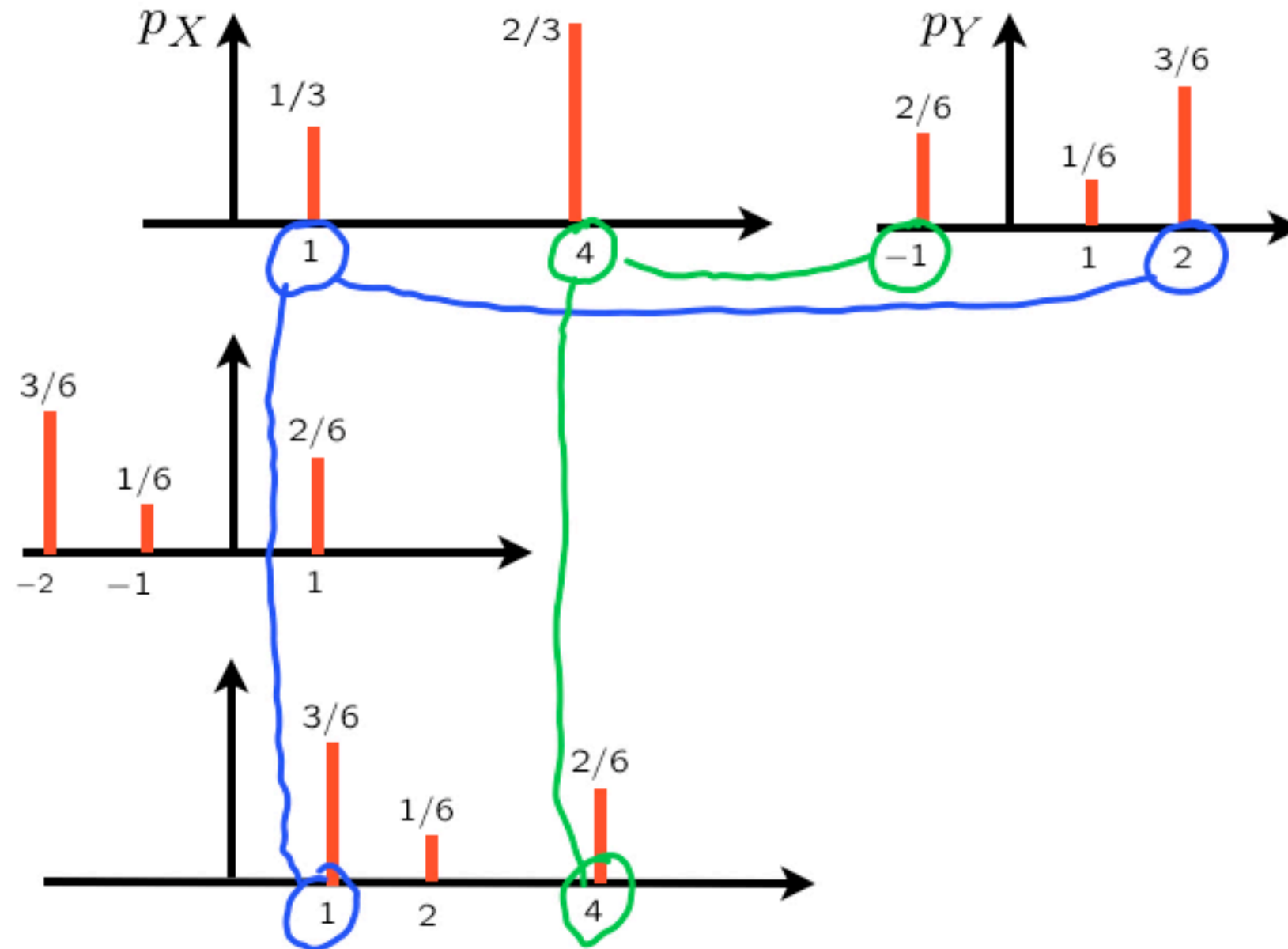
$$p_Z(z) = \sum_x p_X(x) p_Y(z-x)$$



$$p_Z(z) = \sum_x P(X=x, Y=z-x)$$

$$= \sum_x p_X(x) p_Y(z-x)$$

Discrete convolution mechanics



$$p_Z(z) = \sum_x p_X(x) p_Y(z - x)$$

- To find $p_Z(3)$:

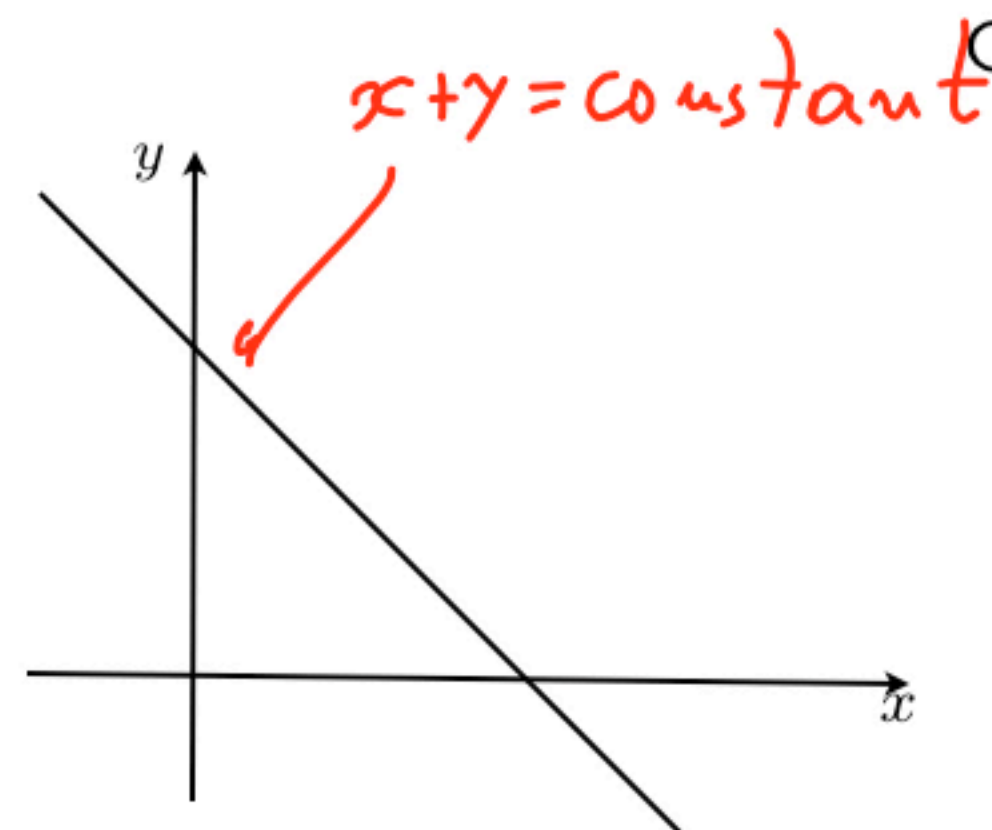
- Flip (horizontally) the PMF of Y
- Put it underneath the PMF of X
- Right-shift the flipped PMF by 3
- Cross-multiply and add
- Repeat for other values of z

The distribution of $X + Y$: the continuous case

- $Z = X + Y$; X, Y independent, continuous
known PDFs

$$p_Z(z) = \sum_x p_X(x) p_Y(z - x)$$

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(\underline{x}) f_Y(\underline{z - x}) dx$$



Conditional on $X = x$: $Z = x + Y$ $x=3$ $Z = Y + 3$

$$f_{Z|X}(z|3) = f_{Y+3|X}(z|3) = f_{Y+3}(z) = f_Y(z-3)$$

$$f_{Z|X}(z|x) = f_Y(z-x)$$

$$f_{X+b}(x) = f_X(x-b)$$

Joint PDF of Z and X :

$$f_{X,Z}(x,z) = f_X(x) f_Y(z-x)$$

From joint to the marginal: $f_Z(z) = \int_{-\infty}^{\infty} f_{X,Z}(x,z) dx$

- Same mechanics as in discrete case (flip, shift, etc.)

The sum of independent normal r.v.'s

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx$$

- $X \sim N(\mu_x, \sigma_x^2)$, $Y \sim N(\mu_y, \sigma_y^2)$, independent

$$Z = X + Y$$

$$f_X(x) = \frac{1}{\sqrt{2\pi} \sigma_x} e^{-(x-\mu_x)^2/2\sigma_x^2} \quad f_Y(y) = \frac{1}{\sqrt{2\pi} \sigma_y} e^{-(y-\mu_y)^2/2\sigma_y^2}$$

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z-x) dx$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi} \sigma_x} \exp\left\{-\frac{(x-\mu_x)^2}{2\sigma_x^2}\right\} \frac{1}{\sqrt{2\pi} \sigma_y} \exp\left\{-\frac{(z-x-\mu_y)^2}{2\sigma_y^2}\right\} dx$$

$$\text{(algebra)} = \frac{1}{\sqrt{2\pi(\sigma_x^2 + \sigma_y^2)}} \exp\left\{-\frac{(z - \mu_x - \mu_y)^2}{2(\sigma_x^2 + \sigma_y^2)}\right\}$$

$$\underbrace{X + Y} + \underbrace{W}_{\cdot}$$

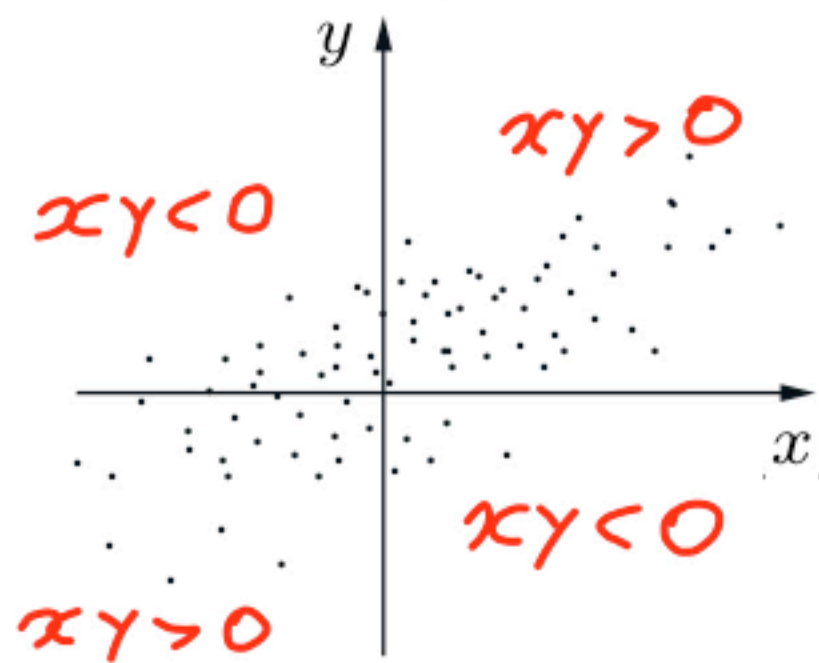
$$N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$$

The sum of finitely many independent normals is normal

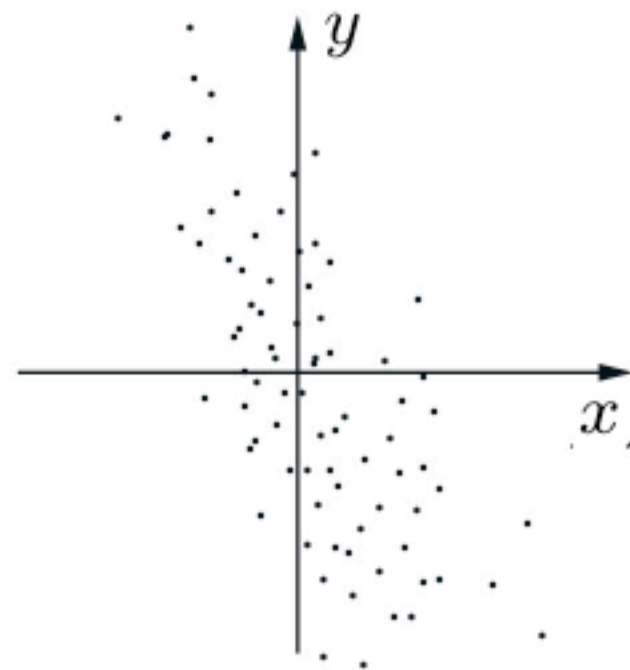
Covariance

- Zero-mean, discrete X and Y
 - if independent: $E[XY] =$

$$= E[X]E[Y] = 0$$



$$E[XY] > 0$$



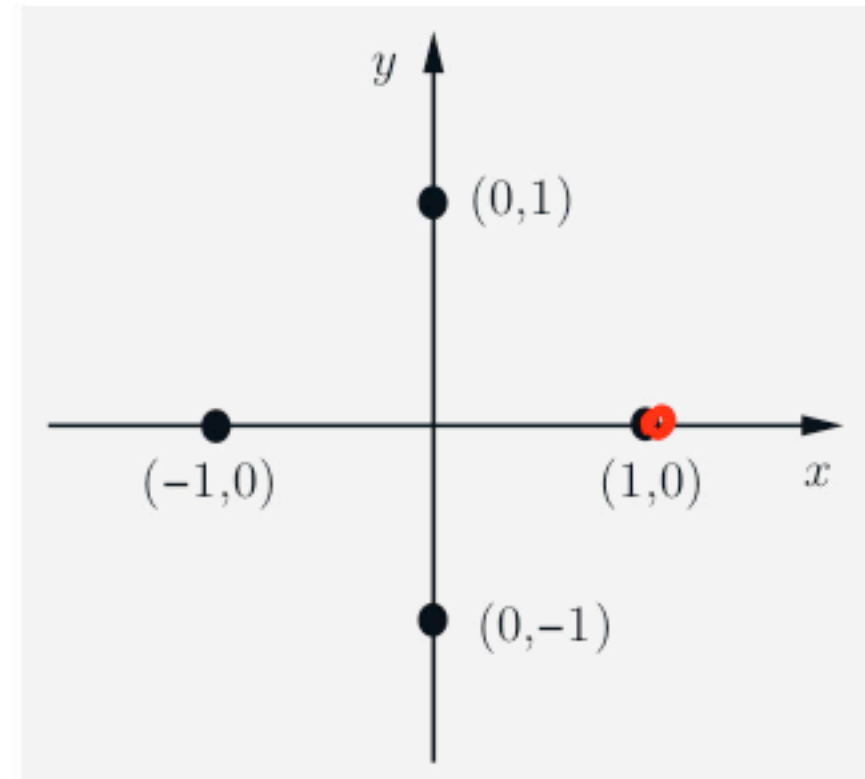
$$E[XY] < 0$$

Definition for general case:

$$\text{cov}(X, Y) = E[(X - E[X]) \cdot (Y - E[Y])]$$

$$\text{and } 0 = E[(X - E[X])] E[Y - E[Y]]$$

- independent $\Rightarrow \text{cov}(X, Y) = 0$
(converse is not true)



$$XY = 0$$

$$\text{Cov} = 0$$

$$X = 1 \Rightarrow Y = 0$$

Covariance properties

$$\begin{aligned}\text{cov}(X, X) &= E[(x - E[x])^2] \\ &= \text{var}(x) = E[x^2] - (E[x])^2\end{aligned}$$

$$\begin{aligned}\text{cov}(aX + b, Y) &= \\ (\text{assume } 0 \text{ means}) \\ &= E[(ax + b)Y] = aE[XY] + bE[Y] \\ &= a \cdot \text{cov}(x, Y)\end{aligned}$$

$$\begin{aligned}\text{cov}(X, Y + Z) &= E[x(y + z)] \\ &= E[xy] + E[xz] = \text{cov}(x, Y) + \text{cov}(x, Z)\end{aligned}$$

$$\text{cov}(X, Y) = E[(X - E[X]) \cdot (Y - E[Y])]$$

$$\begin{aligned}&= E[XY] - E[XE[Y]] \\ &\quad - E[E[X]Y] + E[E[X]E[Y]] \\ &= E[XY] - E[X]E[Y] \\ &\quad - \cancel{E[X]E[Y]} + \cancel{E[X]E[Y]}\end{aligned}$$

$$\text{cov}(X, Y) = E[XY] - E[X]E[Y]$$

The variance of a sum of random variables

$$\begin{aligned}\text{var}(X_1 + X_2) &= E[(X_1 + X_2 - E[X_1 + X_2])^2] \\&= E[(\underbrace{(X_1 - E[X_1])} + \underbrace{(X_2 - E[X_2])})^2] \\&= E[(X_1 - E[X_1])^2 + (X_2 - E[X_2])^2 \\&\quad + 2(X_1 - E[X_1])(X_2 - E[X_2])] \\&= \text{var}(X_1) + \text{var}(X_2) + 2 \text{cov}(X_1, X_2).\end{aligned}$$

The variance of a sum of random variables

$$\text{var}(X_1 + X_2) = \text{var}(X_1) + \text{var}(X_2) + 2 \text{cov}(X_1, X_2)$$

$$\begin{aligned} \text{var}(X_1 + \dots + X_n) &= E[(x_1 + \dots + x_n)^2] \\ (\text{assume 0 means}) &= E\left[\sum_{i=1}^n x_i^2 + \sum_{\substack{i=1, \dots, n \\ j=1, \dots, n \\ i \neq j}} x_i x_j\right] \\ &\quad \left. \vphantom{\sum_{i=1}^n} \right\} n^2 - n \text{ terms} \\ &= \sum_i \text{Var}(x_i) + \sum_{i \neq j} \text{Cov}(x_i, x_j) \end{aligned}$$

$$\text{var}(X_1 + \dots + X_n) = \sum_{i=1}^n \text{var}(X_i) + \sum_{\{(i,j): i \neq j\}} \text{cov}(X_i, X_j)$$

The Correlation coefficient

- Dimensionless version of covariance:

$$-1 \leq \rho \leq 1$$

$$\begin{aligned}\rho(X, Y) &= \mathbf{E} \left[\frac{(X - \mathbf{E}[X])}{\sigma_X} \cdot \frac{(Y - \mathbf{E}[Y])}{\sigma_Y} \right] \\ &= \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}\end{aligned}$$

- Measure of the degree of “association” between X and Y
- Independent $\Rightarrow \rho = 0$, “uncorrelated” (converse is not true)
- $\rho(X, X) = \frac{\text{var}(X)}{\sigma_X^2} = 1$
- $|\rho| = 1 \Leftrightarrow (X - \mathbf{E}[X]) = c(Y - \mathbf{E}[Y])$ (linearly related)
- $\text{cov}(aX + b, Y) = a \cdot \text{cov}(X, Y) \Rightarrow \rho(aX + b, Y) = \frac{a \text{cov}(X, Y)}{|a| \sigma_X \sigma_Y} = \frac{\text{sign}(a)}{\cdot} \rho(X, Y)$

Proof of key properties of the correlation coefficient

$$\rho(X, Y) = \mathbf{E} \left[\frac{(X - \mathbf{E}[X])}{\sigma_X} \cdot \frac{(Y - \mathbf{E}[Y])}{\sigma_Y} \right]$$

$$-1 \leq \rho \leq 1$$

- Assume, for simplicity, zero means and unit variances, so that $\rho(X, Y) = \mathbf{E}[XY]$

$$\begin{aligned} \mathbf{E}[\underline{(X - \rho Y)^2}] &= E[X^2] - 2\rho E[XY] + \rho^2 E[Y^2] \\ 0 \leq &= 1 - 2\rho^2 + \rho^2 = \underline{\underline{1 - \rho^2}} \quad 1 - \rho^2 \geq 0 \Rightarrow \rho^2 \leq 1 \end{aligned}$$

If $|\rho| = 1$, then $X = \rho Y \Rightarrow X = Y$ or $X = -Y$

Interpreting the correlation coefficient

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

- Association does not imply causation or influence

X : math aptitude

Y : musical ability

- Correlation often reflects underlying, common, hidden factor

- Assume, Z , V , W are independent

$$X = \underline{Z} + V \quad Y = \underline{Z} + W$$

Assume, for simplicity, that Z , V , W have zero means, unit variances

$$\text{var}(X) = \text{var}(Z) + \text{var}(V) = 2 \Rightarrow \sigma_X = \sqrt{2} \quad \sigma_Y = \sqrt{2}$$

$$\begin{aligned} \text{cov}(X, Y) &= E[(Z + V)(Z + W)] = E[Z^2] + E[VZ] + E[ZW] + E[VW] \\ &= 1 + 0 + 0 + 0 \end{aligned}$$

$$\rho(X, Y) = \frac{1}{\sqrt{2} \cdot \sqrt{2}} = \frac{1}{2}$$

Correlations matter...

- A real-estate investment company invests \$10M in each of 10 states. At each state i , the return on its investment is a random variable X_i , with mean 1 and standard deviation 1.3 (in millions).

$$\text{var}(X_1 + \dots + X_{10}) = \sum_{i=1}^{10} \text{var}(X_i) + \sum_{\{(i,j): i \neq j\}} \text{cov}(X_i, X_j)$$

$$E[X_1 + \dots + X_{10}] = 10$$

- If the X_i are uncorrelated, then:

$$\text{var}(X_1 + \dots + X_{10}) = 10 \cdot (1.3)^2 = 16.9 \quad \sigma(X_1 + \dots + X_{10}) = 4.1$$

- If for $i \neq j$, $\rho(X_i, X_j) = 0.9$: $\text{cov}(X_i, X_j) = \rho \sigma_{X_i} \sigma_{X_j} = 0.9 \times 1.3 \times 1.3 = 1.52$

$$\text{var}(X_1 + \dots + X_{10}) = 10 \cdot (1.3)^2 + 90 \cdot 1.52 = 154$$

$$\sigma(X_1 + \dots + X_{10}) = 12.4$$