

# Random vectors & common distributions

Sofia Olhede



September 23, 2020

- 1 Random vectors
- 2 Multivariate Random Variables
- 3 Moment Generating Functions

## Random Vectors III

- More generally, we can define the joint frequency/density of random vector formed by a subset of the coordinates of  $X = (X_1 \dots X_d)^T$ , say the first  $k$ ,
- Discrete case:

$$f_{X_1, \dots, X_k}(x_1, \dots, x_k) = \sum_{x_{k+1}} \cdots \sum_{x_d} f_X(x_1, \dots, x_d).$$

- Continuous case:

$$f_{X_1, \dots, X_k}(x_1, \dots, x_k) = \int \cdots \int f_X(x_1, \dots, x_k, y_{k+1}, \dots, y_d) dy_{k+1}, \dots, dy_d.$$

- To marginalize we integrate/sum over the remaining variables from the overall joint density/mass function.
- The  $d$  marginals do not uniquely jointly determine the joint distribution.

# Random Vectors IV

- We may wish to make probabilistic statements about the potential outcomes of one random variable if we already know the outcome of another.
- For this we need the notion of a conditional density/mass function.
- If  $(X_1, \dots, X_d)$  is a continuous/discrete random vector we define the conditional pdf/pmf of  $(X_1, \dots, X_k)$  given  $(X_{k+1} = x_{k+1}, \dots, X_d = x_d)$  as

$$\begin{aligned} f_{X_1, \dots, X_k | X_{k+1}, \dots, X_d}(x_1, \dots, x_k | X_{k+1} = x_{k+1}, \dots, X_d = x_d) \\ = \frac{f_X(x_1, \dots, x_d)}{f_{X_{k+1}, \dots, X_d}(x_{k+1}, \dots, x_d)}, \end{aligned}$$

provided that the denominator is strictly positive.

# Random Vectors V

- The random variables  $X_1, \dots, X_d$  are called independent if and only if for all  $x_1, \dots, x_d$

$$F_{X_1, \dots, X_d}(x_1, \dots, x_d) = F_{X_1}(x_1)F_{X_2}(x_2) \dots F_{X_d}(x_d).$$

- Equivalently the random variables  $X_1, \dots, X_d$  are independent if and only if for all  $x_1, \dots, x_d$

$$f_{X_1, \dots, X_d}(x_1, \dots, x_d) = f_{X_1}(x_1)f_{X_2}(x_2) \dots f_{X_d}(x_d).$$

- For two random variables,  $X$  and  $Y$ , we denote their independence as  $X \perp\!\!\!\perp Y$ .
- Note that when random variables are independent, conditionals reduce to marginals.
- Thus knowing the value of one random variable gives no information on the other.

## Random Vectors VI

- The random vector  $X$  in  $\mathbb{R}^d$  is called conditionally independent of the random vector  $Y$  given the random vector  $Z$  written as

$$X \perp\!\!\!\perp_Z Y \quad \text{or} \quad X \perp\!\!\!\perp Y|Z,$$

if and only if, for all  $x_1, \dots, x_d \in \mathbb{R}$

$$F_{X_1, \dots, X_d|Z, Y}(x_1, \dots, x_d) = F_{X_1, \dots, X_d|Z}(x_1, \dots, x_d) \quad (1)$$

Equivalently this can be reformulated in terms of mass/density functions, as for all  $x_1, \dots, x_d \in \mathbb{R}$

$$f_{X_1, \dots, X_d|Z, Y}(x_1, \dots, x_d) = f_{X_1, \dots, X_d|Z}(x_1, \dots, x_d). \quad (2)$$

- Informally, knowing  $Y$  in addition to  $Z$  provides no additional information about  $X$ . If  $X$  is conditionally independent of  $Y$  given  $Z$  then

$$F_{X, Y|Z} = F_{X|Y, Z} F_{Y|Z} = F_{X|Z} F_{Y|Z}.$$

## Random Vectors VII

- Thus

$$X \perp_Z Y \Leftrightarrow Y \perp_Z X.$$

- Furthermore, if we chose to transform  $\mathbf{X}$  to  $\mathbf{Y}$ , then this can be done from first principles.
- Let  $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be a differentiable bijection.

$$g(\mathbf{x}) = (g_1(\mathbf{x}) \quad \dots \quad g_n(\mathbf{x})), \quad \mathbf{x} = (x_1 \quad \dots \quad x_n)^T \in \mathbb{R}^n.$$

- Let  $X = (X_1 \quad \dots \quad X_n)^T$  have joint density  $f_{\mathbf{X}}(\mathbf{x})$  and define  $\mathbf{Y} = (Y_1 \quad \dots \quad Y_n)^T = g(\mathbf{x})$ . Then with  $\mathcal{Y}^n = g(\mathcal{X}^n)$  and we write the density as

$$f_{\mathbf{Y}}(\mathbf{y}) = f_{\mathbf{X}}(g^{-1}(\mathbf{y})) |\det(J_{g^{-1}}(\mathbf{y}))|, \quad \text{for } \mathbf{y} = (y_1 \quad \dots \quad y_n)^T \in \mathcal{Y}^n,$$

and zero otherwise whenever  $J_{g^{-1}}(\mathbf{y})$  is well-defined.

## Random Vectors VIII

- Here  $J_{g^{-1}}(\mathbf{y})$  is the Jacobian of  $g^{-1}$  i.e. the matrix-valued function

$$J_{g^{-1}}(\mathbf{y}) = \begin{pmatrix} \frac{\partial}{\partial y_1} g_1^{-1}(\mathbf{y}) & \cdots & \frac{\partial}{\partial y_n} g_1^{-1}(\mathbf{y}) \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial y_1} g_n^{-1}(\mathbf{y}) & \cdots & \frac{\partial}{\partial y_n} g_n^{-1}(\mathbf{y}) \end{pmatrix}.$$

- (Sums of random variables). Let  $X$  and  $Y$  be independent continuous random variables with densities  $f_X(x)$  and  $f_Y(y)$  respectively. The density of  $X + Y$  is the convolution of  $f_X(x)$  with  $f_Y(y)$ . Thus

$$f_{X+Y}(u) = \int_{-\infty}^{\infty} f_X(u-v)f_Y(v) dv.$$

- Define  $g : \mathbb{R}^2 \rightarrow \mathbb{R}^2$   $(x, y) \mapsto (x + y, y)$  with inverse transformation  $(u, v) \mapsto (u - v, v)$ . The Jacobian of the inverse is  $\begin{pmatrix} 1 & -1 \\ 0 & 1 \end{pmatrix}$ , with determinant 1.



# Multivariate transformations I

- It follows that

$$f_{X+Y,Y}(u, v) = f_{X,Y}(u - v, v) = f_X(u - v)f_Y(v).$$

We integrate out  $v$  to find the marginal  $f_{X+Y}$ :

$$f_{X+Y}(u) = \int_{-\infty}^{\infty} f_X(u - v)f_Y(v) dv.$$

- The expectation (or expected value) of a random variable  $X$  formalizes the notion of the “average” value taken by that random variable.
- For a continuous random variable this becomes

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} xf_X(x) dx.$$

- For a discrete random variable this becomes

$$\mathbb{E}(X) = \sum_{x \in \mathcal{X}} xf_X(x), \quad \mathcal{X} = \{x \in \mathbb{R} : f_X(x) > 0\}.$$

# Multivariate transformations II

- The expectation satisfies:
- Linearity:  $\mathbb{E}(X_1 + \alpha X_2) = \mathbb{E}(X_1) + \alpha \mathbb{E}(X_2)$ .
- Law of 'unconscious statistician'  $\mathbb{E}(h(X)) = \sum_{x \in \mathcal{X}} h(x) f_X(x)$  (discrete) or  $\mathbb{E}(h(X)) = \int_{x \in \mathcal{X}} h(x) f_X(x)$  (continuous).
- Let  $\mathbf{X} = (X_1 \dots X_d)^T$  be a random vector in  $\mathbb{R}^d$ . For any  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  we define

$$\mathbb{E}\{g(X_1, \dots, X_d)\} = \int_{-\infty}^{\infty} g(x_1, \dots, x_d) f_X(x) dx_1, \dots, dx_d.$$

Similarly in the discrete case

$$\mathbb{E}\{g(X_1, \dots, X_d)\} = \sum_{x_1 \in \mathcal{X}} \dots \sum_{x_d \in \mathcal{X}} g(x_1, \dots, x_d) f_X(x).$$

# Multivariate transformations III

- The mean vector of random vector  $\mathbf{X} = (X_1 \dots X_d)^T$  is defined as

$$\mathbb{E}(\mathbf{X}) = \begin{pmatrix} \mathbb{E}(X_1) \\ \dots \\ \mathbb{E}(X_d) \end{pmatrix},$$

i.e. the vector of means.

- The variance of a random variable  $X$  expresses how disperse the realisations of  $X$  are around its expectation

$$\text{Var}(X) = \mathbb{E}\{(X - \mathbb{E}(X))^2\},$$

if  $\mathbb{E}(X^2)$  is finite.

- Furthermore the covariance of a random variable  $X_1$  with another random variable  $X_2$  expresses the linear dependence between the two. We have

$$\text{Cov}(X_1, X_2) = \mathbb{E}\{(X_1 - \mathbb{E}(X_1))(X_2 - \mathbb{E}(X_2))\}.$$

# Correlation

- The correlation between  $X_1$  and  $X_2$  is defined as

$$\text{corr}(X_1, X_2) = \frac{\mathbb{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1) \text{Var}(X_2)}}.$$

- The correlation conveys equivalent dependence information to the covariance. Advantages: (1) invariant to scale changes, (2) can be understood in absolute terms (ranges in  $[-1, 1]$ ). This is a consequence of the correlation inequality, follows from Cauchy-Schwarz inequality.
- Some useful formulae relating quantities as
  - \*  $\text{Var}(X) = \mathbb{E}(X^2) - \mathbb{E}^2(X) = \mathbb{Cov}(X, X)$ .
  - \*  $\text{Var}(aX + b) = a^2 \text{Var}(X)$ .
  - \*  $\text{Var} \sum_i X_i = \sum_i \text{Var} X_i + \sum_{i \neq j} \mathbb{Cov}(X_i, X_j)$ .
  - \*  $\mathbb{Cov}(X_i, X_j) = \mathbb{E}(X_i X_j) - \mathbb{E}(X_i) \mathbb{E}(X_j)$ .
  - \*  $\mathbb{Cov}(aX_1 + bX_2, Y) = a \mathbb{Cov}(X_1, Y) + b \mathbb{Cov}(X_2, Y)$ .

## Correlation II

- If the second order properties are finite, e.g.  $\mathbb{E}(X_1^2) + \mathbb{E}(X_2^2) < \infty$  then the following are equivalent:
  - \*  $\mathbb{E}(X_1 X_2) = \mathbb{E}(X_1) \mathbb{E}(X_2)$ .
  - \*  $\text{Cov}(X_1, X_2) = 0$ .
  - \*  $\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2)$ .
- Independence implies these three properties. But none of these properties implies independence.
- Let us illustrate this with an example. Let  $X \sim \text{Unif}(-\pi, \pi)$ , and take  $Y = \cos(X)$ . As  $Y$  is a function of  $X$  the two variables cannot be independent.
- They are perfectly dependent, but their covariance is zero.
- We may calculate

$$\Pr(Y > 0) = 1/2, \quad \text{but} \quad [\Pr(Y > 0 | X \in (-\pi, -2))] = 1.$$

## Correlation III

- Despite this we find

$$\mathbb{Cov}(X, Y) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y) \quad (3)$$

$$= \int_{-\pi}^{\pi} x \cos(x) \frac{1}{2\pi} dx - 0 = 0. \quad (4)$$

- Example of how zero correlation does not imply independence.
- Recall that  $\text{Si}(x)$  is the integral whose value is zero at zero of  $\sin(x)/x$  for  $x = 0$ . Let  $X$  and  $Y$  have joint density

$$f_{X,Y}(x, y) = \begin{cases} 1/\pi & \text{if } \text{Si}(x^2 + y^2) \leq 1 \\ 0 & \text{otherwise} \end{cases}.$$

Using symmetry we can directly deduce that  $\mathbb{E}(X) = \mathbb{E}(Y) = 0$ . Thus  $\mathbb{Cov}(X, Y) = \mathbb{E}(XY)$ . But by implementing the integrals we see directly that

$$\mathbb{E}(XY) = 0.$$

# Conditional Expectation I

- We can also calculate the conditional expectation of random variable  $X$  given that of another random variable  $Y$  which took the value  $y$  as

$$\mathbb{E}(X|Y = y) = \begin{cases} \sum_{x \in \mathcal{X}} x \Pr\{X = x|Y = y\} & \text{if } X \text{ and } Y \text{ discrete} \\ \int_{\mathcal{X}} x f_{X|Y=y}(x|y) dx & \text{if } X \text{ and } Y \text{ continuous} \end{cases}$$

- This is the calculation of expectation of the conditional distribution.
- Note that the calculation of  $\mathbb{E}(X|Y = y) = q(y)$  results in a function of  $y$ .
- One can plug  $Y$  into  $q(y)$  and consider  $Y = q(Y)$  as its own random variable.
- Denoted by  $\mathbb{E}(X|Y)$ , this is the formal definition of the conditional expectation.
- Important property/interpretation

$$\mathbb{E}(X|Y) = \arg \min_g \mathbb{E} \|X - g(Y)\|^2.$$

# Conditional Expectation II

- Thus among all measurable functions of  $Y$ ,  $\mathbb{E}(X|Y)$  best approximates  $X$  in the mean square sense.
- Important properties of  $\mathbb{E}(X|Y)$ :
  - \* Unbiasedness  $\mathbb{E}_Y\{\mathbb{E}_{X|Y}(X|Y)\} = \mathbb{E}_X(X)$ .
  - \* If  $X$  is independent of  $Y$  then  $\mathbb{E}(X|Y) = \mathbb{E}(X)$ .
  - \* Taking out known factors:  
 $\mathbb{E}\{g(Y)X|Y\} = g(Y) \mathbb{E}(X|Y)$ .
  - \* Tower property  $\mathbb{E}(\mathbb{E}(X|Y)|g(Y)) = \mathbb{E}(X|g(Y))$ .
  - \* Linearity  $\mathbb{E}(\alpha X_1 + X_2|Y) = \alpha \mathbb{E}(X_1|Y) + \mathbb{E}(X_2|Y)$
  - \* Monotonicity  $X_1 \leq X_2 \Rightarrow \mathbb{E}(X_1|Y) \leq \mathbb{E}(X_2|Y)$ .



## Conditional Expectation III

- The conditional variance of  $X$  given  $Y$  is defined as

$$\mathbb{V}\text{ar}\{X|Y\} = \mathbb{E}_Y\left\{(X - \mathbb{E}_{X|Y}(X|Y))^2|Y\right\} = \mathbb{E}(X^2|Y) - \mathbb{E}^2(X|Y).$$

The law of total variance states that

$$\mathbb{V}\text{ar}(X) = \mathbb{E}_Y(\mathbb{V}\text{ar}(X|Y)) + \mathbb{V}\text{ar}_Y(\mathbb{E}(X|Y)).$$

The proof of this follows directly from

$$\begin{aligned}\mathbb{V}\text{ar}(X) &= \mathbb{E}(X^2) - \mathbb{E}^2(X) \\ &= \mathbb{E}_Y(\mathbb{E}(X^2|Y)) - \mathbb{E}^2(\mathbb{E}(X|Y)) \\ &= \mathbb{E}_Y(\mathbb{V}\text{ar}\{X|Y\} + \mathbb{E}^2(X|Y)) - \mathbb{E}^2(\mathbb{E}(X|Y)) \\ &= \mathbb{E}_Y(\mathbb{V}\text{ar}\{X|Y\}) + \mathbb{E}_Y(\mathbb{E}^2(X|Y)) - \mathbb{E}^2(\mathbb{E}(X|Y)) \\ &= \mathbb{E}_Y(\mathbb{V}\text{ar}(X|Y)) + \mathbb{V}\text{ar}_Y(\mathbb{E}(X|Y)).\end{aligned}$$

QED.

# Conditional Expectation IV

- The covariance matrix of a random vector  $\mathbf{Y} = (Y_1 \dots Y_d)^T$  say  $\mathbf{\Omega} = \{\Omega_{ij}\}$  is a  $d \times d$  symmetric matrix with entries

$$\Omega_{ij} = \text{Cov}\{Y_i, Y_j\} = \mathbb{E}\{(Y_i - \mathbb{E}(Y_i))(Y_j - \mathbb{E}(Y_j))\}, \quad 1 \leq i \leq j \leq d.$$

- Thus it follows that the covariance is the matrix of variance of the variables  $\{Y_i\}$  (on the diagonal), and the covariances of the variables  $\{Y_i\}$  with  $\{Y_j\}$  (on the off-diagonals). We then write

$$\boldsymbol{\mu} = \mathbb{E}\{\mathbf{Y}\} = (\mathbb{E}\{Y_1\} \dots \mathbb{E}\{Y_d\})^T,$$

for the mean vector of  $\mathbf{Y}$ . We also write

$$\text{Var}\{\mathbf{Y}\} = \mathbb{E}\left\{[\mathbf{Y} - \boldsymbol{\mu}][\mathbf{Y} - \boldsymbol{\mu}]^T\right\} = \mathbb{E}\{\mathbf{Y}\mathbf{Y}^T\} - \boldsymbol{\mu}\boldsymbol{\mu}^T.$$

- Thus just like the vector case the expectation of a matrix with random entries is the matrix of expectations of the random entries.

# Covariance calculations

- Let  $\mathbf{Y}$  be a random  $d \times 1$  vector with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Omega}$ .
- For any  $\boldsymbol{\beta} \in \mathbb{R}^d$  we have  $\boldsymbol{\beta}^T \boldsymbol{\Omega} \boldsymbol{\beta} \geq 0$ .
- If  $\mathbf{A}$  is a  $p \times d$  deterministic matrix, then the mean vector and covariance matrix of  $\mathbf{A}\mathbf{Y}$  are  $\mathbf{A}\boldsymbol{\mu}$  and  $\mathbf{A}\boldsymbol{\Omega}\mathbf{A}^T$ , respectively.
- If  $\boldsymbol{\beta} \in \mathbb{R}^d$  is a deterministic vector, then the variance of  $\boldsymbol{\beta}^T \mathbf{Y}$  is  $\boldsymbol{\beta}^T \boldsymbol{\Omega} \boldsymbol{\beta}$ .
- If  $\boldsymbol{\beta}, \boldsymbol{\gamma} \in \mathbb{R}^d$  are deterministic vectors then the covariance of  $\boldsymbol{\beta}^T \mathbf{Y}$  with  $\boldsymbol{\gamma}^T \mathbf{Y}$  is  $\boldsymbol{\gamma}^T \boldsymbol{\Omega} \boldsymbol{\beta}$ .
- Given  $X$  is assumed to be a non-negative random variable. Then, given any  $\epsilon > 0$  we have (Markov's inequality)

$$\Pr(X \geq \epsilon) \leq \frac{\mathbb{E}(X)}{\epsilon}.$$

## Covariance calculations II

- Let  $X$  be a random variable with finite mean  $\mu = \mathbb{E}(X) < \infty$ . Then given any  $\epsilon > 0$  (Chebychev's inequality)

$$\Pr(|X - \mathbb{E}(X)| \geq \epsilon) \leq \frac{\text{Var}(X)}{\epsilon^2}.$$

- For any convex function  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ . If  $\mathbb{E}|\varphi(X)| + \mathbb{E}|X| < \infty$  then Jensen's inequality states

$$\varphi(\mathbb{E}(X)) \leq \mathbb{E}(\varphi(X)).$$

- Let  $X$  be a real random variable with  $\mathbb{E}(X^2) < \infty$ . Let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be a non-decreasing function so that  $\mathbb{E}(g^2(X)) < \infty$ . Then

$$\text{Cov}(X, g(X)) \geq 0.$$

This is a consequence of Chebychev's algebraic inequality.

# Moment Generating Functions

- Let  $X$  be a random variable taking values in  $\mathbb{R}$ . The moment generating function (MGF) of  $X$  is defined as

$$M_X(t) : \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\},$$

and

$$M_X(t) = \mathbb{E}(e^{tX}).$$

- When  $M_X(t)$  and  $M_Y(t)$  exist (and are finite) for  $t \in I$  where  $0 \in I$ . Then

$$* \mathbb{E}|X|^k < \infty \text{ and } \mathbb{E}(X^k) = \frac{d^k M_X}{dt^k}(0) \text{ for all } k \in \mathbb{N}.$$

$$* M_X = M_Y \text{ on } I \text{ if and only if } F_X = F_Y.$$

$$* M_{X+Y}(t) = M_X(t)M_Y(t) \text{ when } X \text{ and } Y \text{ are independent.}$$

- Similarly for a random vector  $\mathbf{X}$  in  $\mathbb{R}^d$  the MGF is

$$M_{\mathbf{X}}(\mathbf{u}) : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}, \quad (5)$$

$$M_{\mathbf{X}}(\mathbf{u}) = \mathbb{E}(e^{\mathbf{u}^T \mathbf{X}}), \quad \mathbf{u} \in \mathbb{R}^d. \quad (6)$$

# Moment Generating Functions II

- A random variable  $X$  is said to follow the Bernoulli distribution with parameter  $p \in (0, 1)$  denoted  $X \sim \text{Bern}(p)$ , if

- \*  $\mathcal{X} = \{0, 1\}$ .

- \*  $f(x; p) = pI(x = 1) + (1 - p)I(x = 0)$ .

The mean, variance and moment generating function of  $X \sim \text{Bern}(p)$  are given by

$$\mathbb{E}(X) = p, \quad \text{Var}(X) = p(1 - p), \quad M_X(t) = 1 - p + pe^t.$$

- A random variable  $X$  is said to follow the Binomial distribution with parameter  $p \in (0, 1)$  and  $n \in \mathbb{N}^+$  denoted  $X \sim \text{Bin}(n, p)$ , if

- \*  $\mathcal{X} = \{0, 1, \dots, n\}$ .

- \*  $f(x; p) = \binom{n}{x} p^x (1 - p)^{n-x}$ .

# Moment Generating Functions III

- The mean, variance and moment generating function of  $X \sim \text{Bin}(n, p)$  are given by

$$\mathbb{E}(X) = np, \quad \text{Var}(X) = np(1 - p), \quad M_X(t) = (1 - p + pe^t)^n.$$

- If  $X = \sum_{i=1}^n Y_i$  where  $Y_i \stackrel{\text{iid}}{\sim} \text{Bern}(p)$  then  $X \sim \text{Bin}(n, p)$ .
- A random variable  $X$  is said to follow the Geometric distribution with parameter  $p \in (0, 1)$  denoted  $X \sim \text{Geom}(p)$ , if
  - \*  $\mathcal{X} = \{0\} \cup \mathbb{N}$ .
  - \*  $f(x; p) = (1 - p)^x p$ .
- The mean, variance and moment generating of  $X \sim \text{Geom}(p)$  are given by

$$\mathbb{E}(X) = \frac{1 - p}{p}, \quad \text{Var}(X) = \frac{1 - p}{p^2}, \quad M_X(t) = \frac{p}{1 - (1 - p)e^t},$$

- the latter for  $t < -\log(1 - p)$ .