

## PROJECT 1 : Some Takeaways

M. PEREIRA

## 1 Random vectors and joint distribution

A **random vector**  $\mathbf{X} = (X_1, \dots, X_n)$  is a vector whose entries  $X_1, \dots, X_n$  are random variables.

Just like we characterize the “randomness” of a random variable using its distribution function, for a random vector we use the **joint distribution function** of its entries, which is the function  $F_{\mathbf{X}} : \mathbb{R}^n \rightarrow [0, 1]$  defined by

$$F_{\mathbf{X}}(x_1, \dots, x_n) = \mathbb{P}(X_1 \leq x_1, \dots, X_n \leq x_n), \quad x_1, \dots, x_n \in \mathbb{R}.$$

In other words,  $F_{\mathbf{X}}(x_1, \dots, x_n)$  is the probability that “simultaneously”  $X_1 \leq x_1$ ,  $X_2 \leq x_2$ , ..., and  $X_n \leq x_n$ .

In some cases, there exists a function  $f_{\mathbf{X}} : \mathbb{R}^n \mapsto \mathbb{R}$  such that  $\forall x_1, \dots, x_n \in \mathbb{R}$ ,

$$F_{\mathbf{X}}(x_1, \dots, x_n) = \int_{-\infty}^{x_n} \dots \int_{-\infty}^{x_1} f_{\mathbf{X}}(y_1, \dots, y_n) dy_1 \dots dy_n, \quad x_1, \dots, x_n \in \mathbb{R},$$

Such a function is called **joint probability density function** (or simply joint density)  $(X_1, \dots, X_n)$ . In particular,  $F_{\mathbf{X}}$  and  $f_{\mathbf{X}}$  are then linked by the relation:

$$f_{\mathbf{X}}(x_1, \dots, x_n) = \frac{\partial^n F}{\partial x_1 \dots \partial x_n}(x_1, \dots, x_n), \quad x_1, \dots, x_n \in \mathbb{R}.$$

**Note:** In the case  $n = 1$ , the joint distribution function and the joint density introduced above correspond to the distribution function and the density of the sole entry of the vector.

## 2 Independence and uncorrelatedness

Two random variables  $X_1$  and  $X_2$  are **independent** if for any values  $x_1, x_2 \in \mathbb{R}$ , the probability that “simultaneously”  $X_1 \leq x_1$  and  $X_2 \leq x_2$  is equal to the product of the probability that  $X_1 \leq x_1$  with the probability that  $X_2 \leq x_2$ :

$$\mathbb{P}(X_1 \leq x_1, X_2 \leq x_2) = \mathbb{P}(X_1 \leq x_1) \cdot \mathbb{P}(X_2 \leq x_2)$$

In this case, the joint distribution function of  $(X_1, X_2)$  is equal to the product of the distribution functions of  $X_1$  and  $X_2$ . If  $X_1$  and  $X_2$  are independent, then the outcome of  $X_1$  has no effect on the outcome of  $X_2$  (and vice-versa).

**Remark:** The notion of independence can be generalized to more than two variables. We say that  $n \geq 2$  random variables  $X_1, \dots, X_n$  are mutually independent if the joint distribution function of  $(X_1, \dots, X_n)$  is equal to the product of the distribution functions of the variables  $X_i$ .

On the other hand, two random variables  $X_1$  and  $X_2$  are called **uncorrelated** if

$$\text{Cov}(X_1, X_2) = \mathbb{E}[X_1 \cdot X_2] - \mathbb{E}[X_1] \cdot \mathbb{E}[X_2] = 0$$

Note that the covariance is just a *measure* of *linear* dependence between two random variables. When faced with two uncorrelated random variables  $X_1, X_2$ , the only thing that we can safely say is that there is no linear dependence between  $X_1$  and  $X_2$ , i.e. that there is no constant  $a \in \mathbb{R}$  such that  $X_2 = aX_1$ . It is for instance perfectly possible that there exists some other non-linear relationship between  $X_1$  and  $X_2$ .

**Example.** Let  $X$  be a random variable following a uniform distribution on  $[-1, 1]$ , meaning that its probability density function is

$$f_X(x) = \begin{cases} 1/2 & \text{if } x \in [-1, 1] \\ 0 & \text{otherwise} \end{cases}$$

Let  $Y$  be the random variable defined by  $Y = X^2$ .

$$\text{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] = \mathbb{E}[X^3] - \mathbb{E}[X]\mathbb{E}[X^2]$$

where

$$\begin{aligned} \mathbb{E}[X] &:= \int_{\mathbb{R}} x f_X(x) dx = \int_{-1}^1 x \cdot \frac{1}{2} dx = \frac{1}{2} \left[ \frac{x^2}{2} \right]_{-1}^1 = 0 \\ \mathbb{E}[X^3] &:= \int_{\mathbb{R}} x^3 f_X(x) dx = \int_{-1}^1 \frac{x^3}{2} dx = \frac{1}{2} \left[ \frac{x^4}{4} \right]_{-1}^1 = 0 \end{aligned}$$

Hence  $\text{Cov}(X, Y) = 0$ , and therefore  $X$  and  $Y$  are uncorrelated. However,  $X$  and  $Y$  are clearly not independent since  $Y$  is a function of  $X$ . The covariance was not able to detect the nonlinear relation between  $X$  and  $Y$ ...

**Advice:** You should keep in mind that independence is a much more stronger assumption than uncorrelatedness. Also:

$$X_1 \text{ and } X_2 \text{ are independent} \Rightarrow X_1 \text{ and } X_2 \text{ are uncorrelated}$$

but in general,

$$X_1 \text{ and } X_2 \text{ are uncorrelated} \not\Rightarrow X_1 \text{ and } X_2 \text{ are independent}$$

### 3 Gaussian variables, vectors and processes

#### 3.1 Gaussian variables and vectors

A **Gaussian variable** (with mean  $\mu$  and variance  $\sigma^2$ ) is a random variable  $X$  whose probability density function  $f_X$  is the function defined by

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad x \in \mathbb{R}$$

In this case, we write  $X \sim \mathcal{N}(\mu, \sigma^2)$ .

If  $X \sim \mathcal{N}(\mu, \sigma^2)$ , then for any constants  $a, b \in \mathbb{R}$ ,  $aX + b \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$ .

Let  $n \geq 2$ . If  $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2), \dots, X_n \sim \mathcal{N}(\mu_n, \sigma_n^2)$  **and**  $X_1, \dots, X_n$  **are (mutually) independent**, then for any constants  $a_1, \dots, a_n \in \mathbb{R}$

$$\sum_{i=1}^n a_i X_i \sim \mathcal{N}\left(\sum_{i=1}^n a_i \mu_i, \sum_{i=1}^n a_i^2 \sigma_i^2\right)$$

A **Gaussian vector** is a random vector  $\mathbf{X} = (X_1, \dots, X_n)$  whose joint probability density function  $f_{\mathbf{X}}$  is the function defined by

$$f_{\mathbf{X}}(y_1, \dots, y_n) = \frac{1}{\sqrt{(2\pi)^n \det \mathbf{\Sigma}}} \exp\left(-\frac{1}{2} \begin{pmatrix} y_1 - \mu_1 \\ \vdots \\ y_n - \mu_n \end{pmatrix}^T \mathbf{\Sigma}^{-1} \begin{pmatrix} y_1 - \mu_1 \\ \vdots \\ y_n - \mu_n \end{pmatrix}\right)$$

where  $\mu_i = \mathbb{E}[X_i]$  ( $1 \leq i \leq n$ ) and  $\mathbf{\Sigma}$  is the covariance matrix of  $\mathbf{X}$ , i.e. the  $n \times n$  matrix defined by whose entry  $(i, j)$  is  $\text{Cov}(X_i, X_j)$  ( $1 \leq i, j \leq n$ ).

For any coefficients  $a_1, \dots, a_n \in \mathbb{R}$ , the random variable defined by

$$\mathbf{X} = (X_1, \dots, X_n) \text{ is a Gaussian vector} \iff \sum_{i=1}^n a_i X_i \quad (1)$$

is a Gaussian variable.

**Remark:** The entries of a Gaussian vector are Gaussian variables (just take  $a_i = 1$  and setting all the other coefficients to 0 in (1)). However, the fact that a vector  $\mathbf{X}$  is composed of entries which are Gaussian variable is in general NOT enough to conclude that  $\mathbf{X}$  is Gaussian vector! You must also have that (1) is satisfied for any coefficients!

#### TIP

To show that a random vector  $\mathbf{X}$  is a Gaussian vector, you can show that for any choice of coefficients  $a_1, \dots, a_n$  the random variable (1) is a Gaussian variable.

If  $\mathbf{X} = (X_1, \dots, X_n)$  is a Gaussian vector, then for any entries  $X_i$  and  $X_j$ ,

$$\text{Cov}(X_i, X_j) = 0 \Rightarrow X_i \text{ and } X_j \text{ are independent.}$$

#### TIP

To show that two random variables  $X_1$  and  $X_2$  are independent, you can show that the vector  $(X_1, X_2)$  is a Gaussian vector and that  $\text{Cov}(X_1, X_2) = 0$ .

**Note:** To use this tip, it is not enough to show that  $X_1$  and  $X_2$  are Gaussian variables! Indeed, two uncorrelated Gaussian random variables are not necessarily independent<sup>1</sup>. But if, besides, they form a Gaussian vector, then they will be independent.

### 3.2 Gaussian processes

A **Gaussian process** is a process  $(X_t, t \in \mathbb{Z})$  such that for any  $n \geq 1$  times  $t_1, \dots, t_n \in \mathbb{Z}$  the vector  $(X_{t_1}, \dots, X_{t_n})$  is a Gaussian vector.

**Note:** For a process  $(X_t, t \in \mathbb{Z})$  to be a Gaussian process it is NOT enough to only check that for any  $t \in \mathbb{Z}$ ,  $X_t$  is a Gaussian variable.

$$\begin{aligned}
 & \text{For any number } n \geq 1 \text{ of arbitrary times } t_1, \dots, t_n \in \mathbb{Z} \text{ and any choice of coefficients } a_1, \dots, a_n \in \mathbb{R}, \text{ the random variable} \\
 & \sum_{i=1}^n a_i X_{t_i} \quad (2) \\
 & \text{is a Gaussian variable.}
 \end{aligned}$$

$(X_t, t \in \mathbb{Z})$  is a Gaussian process  $\iff$

#### TIP

To show that  $(X_t, t \in \mathbb{Z})$  is a Gaussian process, you have to show that for any  $n \geq 1$ , any times  $t_1, \dots, t_n \in \mathbb{Z}$ , and any coefficients  $a_1, \dots, a_n \in \mathbb{R}$ , the random variable (2) is a Gaussian variable.

If  $(X_t, t \in \mathbb{Z})$  is a Gaussian process, then for any times  $t_1, t_2 \in \mathbb{Z}$ ,

$$\text{Cov}(X_{t_1}, X_{t_2}) = 0 \Rightarrow X_{t_1} \text{ and } X_{t_2} \text{ are independent.}$$

<sup>1</sup>cf. [https://en.wikipedia.org/wiki/Normally\\_distributed\\_and\\_uncorrelated\\_does\\_not\\_imply\\_independent](https://en.wikipedia.org/wiki/Normally_distributed_and_uncorrelated_does_not_imply_independent)