



Introduction to statistics and data science

A hands-on approach

Stéphane Puechmorel [†]

October 8, 2025

[†] ENAC Dept. SINA

email: stephane.puechmorel@enac.fr

phone: 0562259503

Introduction

WORK ORGANIZATION

- All lectures will include course material and exercises.
- The notions are illustrated using real data and statistical software.
- Should you encounter any difficulties, please do not hesitate to inquire.
- Personal work is essential and must be completed regularly.

Visual exploration of data

VECTORS

Operators

- Operators and functions are applied elementwise.

```
> a <- 1:10
> b <- 1:10
> a+b
[1] 2 4 6 8 10 12 14 16 18 20
> a*b
[1] 1 4 9 16 25 36 49 64 81 100
> sin(a)
[1] 0.8414710 0.9092974 0.1411200 -0.7568025 -0.9589243 -0.2794155
[7] 0.6569866 0.9893582 0.4121185 -0.5440211
```

Figure 8: Operating on vectors.

- This is true also for comparison operators.

```
> c(1,3,5,8) <= c(0,4,5,8)
[1] FALSE TRUE TRUE TRUE
```

Figure 9: Comparison operators.

VECTORS

Accessing elements

- An element in a vector can be referred to by its index.

```
> a<-1:10
> a[1]
[1] 1
> a[2:4]
[1] 2 3 4
> a[c(2,6,8)]
[1] 2 6 8
```

Figure 10: Accessing elements.

- A selection by booleans is also possible.

```
> a[c(TRUE, FALSE, TRUE)]
[1] 1 3 4 6 7 9 10
> a[a >= 5]
[1] 5 6 7 8 9 10
```

Figure 11: Boolean selection.

DATA FRAMES

A convenient object.

- Data frames are arrays holding observed values.
- Observed characteristics are in columns, observations in rows.
- The columns are named and can be referred to by their names.

```
> dataset <- data.frame(id=1:100, val=0:0.01:99)
> str(dataset["id"])
'data.frame':   100 obs. of  1 variable:
 $ id: int   1  2  3  4  5  6  7  8  9 10 ...
> str(dataset$val)
int [1:100] 0 1 2 3 4 5 6 7 8 9 ...
```

Figure 12: A data frame.

- Data frames are used by PowerBI to communicate with R.

PASSENGERS DATASET

Data collection

- The dataset, collected from the site <https://data.europa.eu/data/datasets?locale=en> is stored in the file `avia_paoc_monthly.csv` that can be retrieved on ecampus in the folder `Datasets`.
- Launch PowerBI and select Get data from other sources.
- Select Text/CSV and open the file.
- A table view of the dataset appears.

PASSENGERS DATASET

Measures

- A measure is a value that can be computed from the dataset.
- It may be viewed as a practical implementation of a random variable, a concept introduced later in the course.
- Since measures are computed on-the-fly, they always reflect the most up-to-Date dataset.
- Using them is a must in PowerBI.
- Create a `New quick measure` on the "passengers" dataset that sums the amount of carried passengers by quarter for a given country (select the one you prefer).
- Use a visual to plot the result.
- Custom measures can be used, provided you know the "DAX" language.

POWERBI AND R

R as a data processor

- Open the `Get data` tab, then select `Run R script`.
- In the script box, type:

```
1 output <- dataset
2 output$z <- as.numeric(dataset$y)+1
```

- A new column named z was added and its values are those of column y plus 1.

Probability theory

A BIT OF HISTORY

A modern view

- In 1933, the Russian mathematician Andrey Kolmogorov gives the first axiomatic description of probability theory.
- This is the approach we will take in the sequel.

Challenges

- To fully understand Kolmogorov's contribution, we must consider cases in which counting the number of favorable outcomes is insufficient.
- For example, take a situation where the measure of interest is a real number, such as a current or a pressure.
- The value itself is not discrete; only the number of times it falls within a given interval can be observed.
- The Fermat-Pascal approach must be extended to these cases.

EXERCISE

Two factories, denoted F_0, F_1 are producing components. The end customer receives a component that comes from F_0 with probability 0.3. Furthermore, the probability of the component being defective is 0.1 when coming from factory F_0 and 0.06 when coming from factory F_1 . If the received component is defective, what is the probability that it has been produced by F_1 ?

EXERCISE

Bernoulli random variable

- A bernoulli random variable X can take only two values: 0 and 1.
- If $P(X = 1) = p$, can you compute $P(X = 0)$?
- Compute $E[X]$.

Binomial random variable

- A binomial random variable X with parameters $n \in \mathbb{N}, p \in [0, 1]$ has distribution:

$$p(k) = \binom{n}{k} p^k (1-p)^{n-k}, k = 0 \dots n. \quad (11)$$

- Find $E[X]$.

MOMENTS

The expected value in a general sense

- If $g: \mathbb{R} \rightarrow \mathbb{R}$ is piecewise continuous and

$$\sum_{x \in E} |g(x)| p(x) < +\infty, \quad (14)$$

then one defines:

$$E[g(X)] = \sum_{x \in E} g(x) p(x) \quad (15)$$

- The variance of the discrete random variable X is:

$$V(X) = E \left[(X - E[X])^2 \right] \quad (16)$$

- It measures the dispersion of X around its expected value.

MOMENTS

Higher order moments

- Let n be an integer. The n -th moment of X is, if it exists:

$$E[X^n]. \quad (17)$$

- The moment generating function is the function:

$$t \mapsto m_X(t) = E[e^{tX}] \quad (18)$$

- If its domain contains 0, then:

$$E[X^n] = m^{(n)}(0), \quad (19)$$

where $m^{(n)}$ denotes the n -th derivative of m_X .

CONDITIONING

Independence

- By the Bayes' formula:

$$P(X = x|Y = y) = \frac{P(X = x, Y = y)}{P(Y = y)}. \quad (20)$$

- If $P(X = x, Y = y) = P(X = x)P(Y = y)$, then $P(X = x|Y = y) = P(X = x)$.
- In such a case, X, Y are said to be **independent**.
- If X, Y are independent, then $E[XY] = E[X]E[Y]$.
- If X_1, \dots, X_n are pairwise independent, then:

$$V(X_1 + \dots + X_n) = V(X_1) + \dots + V(X_n). \quad (21)$$

EXERCISE

A pair of dice... again !

- Let X, Y be two independent random variables corresponding to the respective values of two die throws. Compute :

$$E[X], E[Y], E[X + Y].$$

- Same question with the variances.

A taste of estimation theory

- Let X_1, \dots, X_n be a sequence of independent random variables and let:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

- Compute $E[\bar{X}], V(\bar{X})$. What happens if $n \rightarrow +\infty$?

POISSON DISTRIBUTION

A model for random events

- Given a time interval of length $[t_0, t_1]$, one counts the number of occurrences X of an event.
- The random variable X is said to have a Poisson distribution with rate λ if:

$$P(X = k) = \frac{(\lambda T)^k}{k!} e^{-\lambda T}, \quad (22)$$

where $T = t_1 - t_0$.

- $E[X] = V(X) = \lambda T$.
- The Poisson distribution models a situation where the occurrences of an event are independent and occur at constant rate.

EXERCISE

Distribution of a sum

- The moment generating function (MGF) of a random variable is the function $m_X(t) = E[e^{tX}]$. It is characteristic of a distribution. Prove that if X, Y are independent, $m_{X+Y}(t) = m_X(t)m_Y(t)$.
- Compute the MGF of X with Poisson distribution of rate λ .
- Let X, Y be independent random variables with Poisson distributions of respective rates λ, μ . Show that $X + Y$ has Poisson distribution of rate $\lambda + \mu$.

Selecting events

Let X be a random variable with Poisson distribution of rate λ . The underlying events are selected at random with probability p and counted, yielding a random variable Y . Prove that Y has Poisson distribution with rate $p\lambda$.

REAL RANDOM VARIABLES

Definition

- Let (Ω, \mathcal{T}, P) be a measure space. A mapping $X: \Omega \rightarrow \mathbb{R}$ is said to be a real random variable if, for any $t \in \mathbb{R}$, $\{\omega | X(\omega) \leq t\}$ is an event in \mathcal{T} .
- If X is a real random variable, its cumulative distribution function (CDF) is the mapping:

$$F_x: t \in \mathbb{R} \mapsto P(X \leq t). \quad (23)$$

- $\lim_{t \rightarrow -\infty} F_x(t) = 0, \lim_{t \rightarrow +\infty} F_x(t) = 1.$
- F_x is right continuous and admits a limit to the left.
- If F_x is differentiable, its derivative is called the density of X .

REAL RANDOM VARIABLES

Expected value

- Let X be a real random variable with density p_X .
- If $g: \mathbb{R} \rightarrow \mathbb{R}$ is piecewise continuous, then one defines:

$$E[g(X)] = \int_{\mathbb{R}} g(t)p_X(t)dt \quad (24)$$

provided:

$$\int_{\mathbb{R}} |g(t)|p_X(t)dt < +\infty. \quad (25)$$

- As special cases:

$$P(X \in [a, b]) = E[1_{[a, b]}(X)] = \int_a^b p_X(t)dt \quad (26)$$
$$F_X(t) = \int_{-\infty}^t p_X(t)dt.$$

REAL RANDOM VARIABLES

Joint CDF and density

- If X, Y are real random variables, the CDF of the couple (X, Y) is the mapping:

$$F_{X,Y}: (s, t) \in \mathbb{R}^2 \mapsto P(X \leq s, Y \leq t). \quad (27)$$

- Taking the derivative with respect to s, t allows defining the joint density:

$$p_{X,Y}(x, y) = \frac{\partial^2}{\partial_s \partial_t} F_{X,Y}(t, s)|_{x,y}. \quad (28)$$

- X, Y are said to be independent if:

$$F_{X,Y}(s, t) = F_X(s)F_Y(t). \quad (29)$$

MARGINAL DENSITIES

- Let X, Y be a couple of real random variables with joint density $p_{X,Y}$.
- The density of X (resp. Y) can be obtained as:

$$p_X(x) = \int_{\mathbb{R}} p_{X,Y}(x, t) dt \text{ (resp.) } p_Y(y) = \int_{\mathbb{R}} p_{X,Y}(s, y) ds \quad (30)$$

- The conditional density of X knowing Y is defined as:

$$p(X|Y = y)(s) = \frac{p_{X,Y}(s, y)}{p_Y(y)}. \quad (31)$$

- If X, Y are independent, $p(X|Y = y)(s) = p_X(s)$ or, equivalently, $p_{X,Y}(x, y) = p_X(x)p_Y(y)$.

CHARACTERISTIC FUNCTION

Definition

- Let X be a real random variable. Its characteristic function is:

$$\phi_X: t \in \mathbb{R} \mapsto E[e^{itX}]. \quad (32)$$

- Using the density p_X , it can be computed as:

$$\phi_X(t) = \int_{\mathbb{R}} e^{itx} p_X(x) dx. \quad (33)$$

Properties

- If two random variables have the same characteristic function, they have the same distribution.
- If X, Y are independent, $\phi_{X+Y}(t) = \phi_X(t)\phi_Y(t)$.
- $E[X^k] = i^{-k} \phi_X^{(k)}(0), k \in \mathbb{N}$.

UNIFORM DISTRIBUTION

- The uniform distribution on an interval $[a, b]$ has density:

$$p(x) = \begin{cases} \frac{1}{b-a} & x \in [a, b], \\ 0 & \text{otherwise.} \end{cases} \quad (34)$$

- If X has uniform distribution on $[a, b]$, then:

$$\begin{aligned} E[X] &= \frac{a+b}{2}, \quad V(X) = \frac{(b-a)^2}{12}, \\ \phi_X(t) &= \frac{e^{itb} - e^{ita}}{t(b-a)}. \end{aligned} \quad (35)$$

NORMAL DISTRIBUTION

Density

- Let $\mu \in \mathbb{R}, \sigma > 0$. A real random variable X is said to have a normal distribution $\mathcal{N}(\mu, \sigma^2)$ if its density is:

$$p_X(t) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(t-\mu)^2}{2\sigma^2}}. \quad (36)$$

- If $X \sim \mathcal{N}(\mu, \sigma^2)$, then:

$$Y = \frac{X - \mu}{\sigma} \sim \mathcal{N}(0, 1). \quad (37)$$

NORMAL DISTRIBUTION

Moments

If $X \sim \mathcal{N}(\mu, \sigma^2)$:

$$\begin{aligned} E[X] &= \mu, \quad V(X) = \sigma^2, \\ \phi_X(t) &= \exp(i\mu t - \sigma^2 t^2 / 2). \end{aligned} \tag{38}$$

Cumulative distribution function

- There is no closed-form expression for the CDF of a $\mathcal{N}(\mu, \sigma^2)$ distribution, but it can be evaluated numerically.
- If Φ is the CDF of a $\mathcal{N}(0, 1)$ distribution, then the CDF of a $\mathcal{N}(\mu, \sigma^2)$ distribution is:

$$\Phi\left(\frac{x - \mu}{\sigma}\right). \tag{39}$$

CENTRAL LIMIT THEOREM

Convergence in distribution

A sequence $(X_n)_{n \in \mathbb{N}}$ of real random variables is said to converge in distribution to a random variable X if:

$$\forall t \in \mathbb{R}, \lim_{n \rightarrow +\infty} P(X_n \leq t) = P(X \leq t). \quad (40)$$

Theorem (Central limit)

If $(X_n)_{n \in \mathbb{N}}$ is a sequence of independent, identically distributed real random variables, then the sequence of random variables $\sqrt{n}(\bar{X}_n - \mu)$ converges in distribution to $X \sim \mathcal{N}(0, \sigma^2)$, with:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i, \mu = E[X_i], \sigma^2 = V(X_i). \quad (41)$$

EXPONENTIAL DISTRIBUTION

- The exponential distribution with rate $\lambda > 0$ has density:

$$p(x; \lambda) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0, \\ 0, & x < 0. \end{cases} \quad (42)$$

- If X is exponentially distributed with rate λ :

$$E[X] = \frac{1}{\lambda}, V(X) = \frac{1}{\lambda^2}, \phi_X(t) = \frac{\lambda}{\lambda - it}. \quad (43)$$

- The exponential distribution is memoryless:

$$P(X > t + s | X > s) = P(X > t), \quad s, t \geq 0. \quad (44)$$

VECTOR RANDOM VARIABLES

CDF and density

- Let (Ω, \mathcal{T}, P) be a measure space. Let $X: \Omega \rightarrow \mathbb{R}^n$ be a mapping. X is said to be a vector (or multivariate) random variable iff:

$$\{\omega | X_1(\omega) \leq t_1, \dots, X_n(\omega) \leq t_n\} \in \mathcal{T}, (t_1, \dots, t_n) \in \mathbb{R}^n. \quad (45)$$

- The CDF of X is:

$$F_x(t_1, \dots, t_n) = P(X_1 \leq t_1, \dots, X_n \leq t_n). \quad (46)$$

- Taking the partial derivatives yields the density:

$$p_X(t_1, \dots, t_n) = \frac{\partial^n F_X(t_1, \dots, t_n)}{\partial t_1 \dots \partial t_n}. \quad (47)$$

Estimation

INTRODUCTORY EXAMPLE

Modeling the delays

-