# Statistics for Data Science

Marina Anca Cidota
Faculty of Mathematics and Computer Science
University of Bucharest

✉ cidota@fmi.unibuc.ro

# References

1. Robert Shumway, David Stoffer. Time Series Analysis and its Applications with R Examples, 4th edition, Springer, 2017
2. Peter Brockwell, Richard Davis. Time Series: Theory and Methods, Springer-Verlag, 1987
3. James Hamilton. Time Series Analysis, Princeton University Press, 1994
4. Norman Matloff. Probability and Statistics for Data Science, CRC Press, 2019.
5. Monica Dumitrecu, Anton Batatorescu. Applied Statistics using the R system, Ed. Universitatii din Bucuresti, 2006.
6. Ion Vaduva. Analiza dispersionala, Ed. Tehnica, 1970
7. Peter Bruce, Andrew Bruce. Practical Statistics for Data Scientists, O'Reilly, 2017
8. Hadley Wickham, Garrett Grolemund. R for Data Science, O'Reilly, 2016.

# Introduction

Def. $(\Omega, K, P)$ is called a probability space, where

1) $\Omega$ — sample space = the set of all possible outcomes

2) $K \subset P(\Omega)$ the power set of $\Omega$ (the set of all subsets of $\Omega$)

$K$ is a $\sigma$-algebra
$$\begin{cases} \emptyset \in K \\ A \in K \Rightarrow \overline{A} \in K \\ A_n \in K \Rightarrow \bigcup_{n=1}^{\infty} A_n \in K \end{cases}$$

3) $P : K \longrightarrow [0,1]$ is a probability measure function

$P(\Omega) = 1$

— a measure function: $\begin{cases} P : K \longrightarrow [0, \infty) \\ P(\emptyset) = 0 \\ P\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n) \text{ if } A_n \in K, \\ \qquad\qquad A_i \cap A_j = \emptyset, \ \forall \, i \neq j. \end{cases}$

The elements of $K$ are called random events.

3

# Introduction

Def. if $S$ is a set, the $\sigma$-algebra generated by $S$, written $B(S)$, is:

$$B(S) \stackrel{def}{=} \bigcap_{\substack{\mathcal{F} \supset S \\ \mathcal{F} \text{ is } \sigma\text{-algebra}}} \mathcal{F}$$

If $S = \mathbb{R}$ then $\mathcal{S} \stackrel{not}{=} B(S) = B\left((a,b) \mid a,b \in \mathbb{R}\right)$

If $S$ is finite or countable then $\mathcal{S} = P(S)$.

$(S, \mathcal{S})$ is a measurable space.

Def. $X: (\Omega, \mathcal{X}) \rightarrow (S, \mathcal{S})$ is a random variable if

$\forall B \in \mathcal{S} \implies X^{-1}(B) \in \mathcal{X}$.

$P \circ X^{-1}: \mathcal{S} \rightarrow [0,1]$, $\left(P \circ X^{-1}\right)(B) = P\left(X^{-1}(B)\right)$

4

# Introduction

The cumulative distribution function (CDF) of $X$ is

$$F(x) \overset{def}{=} (P \circ X^{-1})(-\infty, x) = P(X < x)$$

$$F : \mathbb{R} \to [0, 1]$$

$$F(-\infty) = 0; \quad F(\infty) = 1; \quad F \text{ is monotonically increasing}$$

- for $X$ a continuous random variable, $P \circ X^{-1}$ is specified by the density function $f(x)$,

$$f(x) \geq 0 \quad \forall x$$

$$\int_{-\infty}^{\infty} f(x) \, dx = 1$$

$$f(x) = F'(x)$$

$$F(x) = \int_{-\infty}^{x} f(t) \, dt$$

# Introduction

- for $X$ a discrete random variable, $P \circ X^{-1}$ is specified by $\{p(x), x \in S\}$, $p(x) \geq 0$, $\forall x$

$$\sum_{x \in S} p(x) = 1$$

$$X = \begin{pmatrix} x_1 & \cdots & x_n \\ p_1 & \cdots & p_n \end{pmatrix}$$

$$F(x) = \sum_{x_i < x} P(x = x_i)$$

$$E(X) = \int_{-\infty}^{\infty} x \, f(x) \, dx \qquad \left( \text{or } \sum_{i=1}^{n} x_i p_i \text{ for the discrete case} \right)$$

$$Var(x) = E\left( (x - E(x))^2 \right) = E(x^2) - E(x)^2$$

$$T = \sqrt{Var(X)} \quad - \text{ standard deviation}$$

# Stochastic models of different phenomena of interest

- Categorical random variables/vectors.

- Quantitative discrete/continuous random variables/vectors.

Def. A stochastic process is a family of random variables $\{X_t, t \in T\}$ defined on a probability space $(\Omega, K, P)$.

time $T$, states $S$

$\left\langle \begin{array}{l} \text{time } \mathbb{N}/\mathbb{Z}/\mathbb{R}/[0,\infty) \\ \text{state discrete}/\mathbb{R}/\mathbb{R}^d \end{array} \right.$

# Stochastic models of different phenomena of interest

- Observations in "transversal" studies (at a fixed moment of time)

  $\{X_1, \dots, X_n\}$ random variables, independent and identically distributed $(iid)$ like the stochastic model $X: \Omega \to S$. Statistical data are the observed values $(x_1, \dots, x_n) \in S^n$.

- Observations in "longitudinal" studies (on a fixed time interval)

  $\{X_t, t \leq t_n\} \subset \{X_t, t \in T\}$

  statistical data are the observed part of the process trajectory $(x_t, t \leq t_n)$

# Stochastic models of different phenomena of interest

Parametric statistics

Hypothesis: statistical data come from a stochastic model whose distribution has a known functional form, but it depends on an unknown parameter $\theta \in \mathbb{R}^k$, $k \geq 1$.

- Statistical data are available;
- Looking for estimators of $\theta$ or tests on hypotheses over $\theta$.

# Example of stochastic models

1. categorical 1-dim random variable

   $X$ = satisfaction degree

   $S = \{vun, un, sa, vsa\}$

   $$P \circ X^{-1} = \begin{pmatrix} vun & un & sa & vsa \\ p_1 & p_2 & p_3 & p_4 \end{pmatrix} \qquad \begin{array}{l} p_i \geq 0 \\ \sum\limits_{i=1}^{4} p_i = 1 \end{array}$$

2. categorical 2-dim random vector

   $X = (X_1, X_2)' = ($ intention to vote, satisfaction degree about the current situation$)$

   $S = \{$ (yes, vun), (yes, un), (yes, sa), (yes, vsa),
   
   (no, vun), (no, un), (no, sa), (no, vsa) $\}$

# Example of stochastic models

Quantitative discrete random variables

3. discrete uniform distribution — e.g. rolling a dice

$$X \sim U\{1,\ldots,r\}, \quad r \in N, \quad r \geq 2$$

$$S = \{1, 2, \ldots, r\}$$

$$P(X = x) = \frac{1}{r}, \quad x = 1, 2, \ldots, r$$

$$E(X) = \frac{r+1}{2}$$

$$Var(X) = \frac{r^2 - 1}{12}$$

# Example of stochastic models

4. Bernoulli distribution — e.g. having "success" in a trial with two possible outcomes: success/failure.

$$X \sim B(1, \theta), \quad \theta \in (0,1) \quad S = \{0,1\}$$

$$P_\theta(X=x) = \theta^x (1-\theta)^{1-x}, \quad x=0,1$$

$$X : \begin{pmatrix} 0 & 1 \\ 1-\theta & \theta \end{pmatrix}$$

$$E(X) = \theta$$

$$Var(X) = \theta(1-\theta)$$

# Example of stochastic models

5. Binomial distribution — e.g. number of "successes" in $n$ independent trials with two possible outcomes: success/failure.

$$X \sim Bi(n, \theta), \quad \theta \in (0, 1), \quad S = \{0, 1, \ldots, n\}$$

$$P_\theta(X = x) = C_n^x \cdot \theta^x (1-\theta)^{n-x}, \quad x = 0, 1, \ldots, n$$

$$E(X) = n \cdot \theta$$

$$Var(X) = n \cdot \theta(1-\theta)$$

# Example of stochastic models

6. Poisson distribution – e.g. number of defective devices in a very large volume lot.

$$X \sim Po(\theta), \quad \theta \in (0, \infty), \quad S = \mathbb{N}$$

$$P_\theta(X = x) = \frac{\theta^x}{x!} e^{-\theta}, \quad x = 0, 1, \ldots$$

$$E(X) = \theta$$

$$Var(X) = \theta$$

# Example of stochastic models

7. Geometric distribution — e.g. the moment of the first "success" in a sequence of independent trials with two possible outcomes: success / failure.

$$X \sim Ge(\theta), \quad \theta \in (0,1), \quad S = \mathbb{N}^*$$

$$P_\theta(X = x) = (1-\theta)^{x-1} \cdot \theta \qquad x = 1, 2, \ldots$$

$$E(X) = \frac{1}{\theta}$$

$$Var(X) = \frac{1-\theta}{\theta^2}$$

# Example of stochastic models

Quantitative discrete random vectors

8. Multinomial distribution - e.g. sampling with replacement

$n$ trials - extractions from a ballot box containing balls of $d$ colors

$$X = (X_1, \ldots, X_d)' \sim M(n; p_1, \ldots, p_d)$$

$$P_\theta\left(X = (x_1, \ldots, x_d)'\right) = \frac{n!}{x_1! \cdots x_d!} \, p_1^{x_1} \cdots p_d^{x_d}, \text{ where}$$

$$S = \left\{ (x_1, \ldots, x_d)' \mid x_i \in \{0, 1, \ldots, n\} \; \forall i = \overline{1,d}, \; \sum_{i=1}^{d} x_i = n \right\}$$

$$\theta = (p_1, \ldots, p_d), \quad p_i \in [0,1] \; \forall i = \overline{1,d}, \; \sum_{i=1}^{d} p_i = 1$$

$$E(X_i) = n p_i$$
$$i = \overline{1,d}$$
$$Var(X_i) = n p_i (1 - p_i) \qquad\qquad Cov(X_i, X_j) = - n p_i p_j, \quad i \neq j$$

# Example of stochastic models

Quantitative continuous random variables

9. Continuous uniform distribution - e.g throw a dart at random at the interval $[0, \theta]$ - all the points are equally like to be hit.

$$X \sim U(0, \theta), \quad \theta \in (0, \infty)$$

$$f(x; \theta) = \begin{cases} \frac{1}{\theta} & , \ x \in [0, \theta] \\ 0 & , \ x \notin [0, \theta] \end{cases}$$

$$E(x) = \frac{\theta}{2}$$

$$Var(x) = \frac{\theta^2}{12}$$

17

# Example of stochastic models

10. Exponential distribution — e.g. the simplest model for lifetime

$$X \sim \text{Expo}(\theta), \quad \theta \in (0, \infty)$$

$$f(x; \theta) = \begin{cases} \frac{1}{\theta} \exp\left(-\frac{x}{\theta}\right), & x \in [0, \infty) \\ 0, & x \in (-\infty, 0) \end{cases}$$

$$E(X) = \theta$$

$$\text{Var}(X) = \theta^2$$

# Example of stochastic models

11. Gamma distribution — e.g. general model for lifetime

$$X \sim \text{Gamma}(\alpha, \theta), \quad \alpha \in (0, \infty), \quad \theta \in (0, \infty)$$

$$f(x; \alpha, \theta) = \begin{cases} \frac{1}{\Gamma(\alpha) \cdot \theta^{\alpha}} \cdot x^{\alpha-1} \cdot \exp\left(-\frac{x}{\theta}\right), & x \in \langle 0, \infty) \\ 0 & , \quad x \in (-\infty, 0) \end{cases}, \text{ where}$$

$$\Gamma(\alpha) = \int_{0}^{\infty} x^{\alpha-1} e^{-x} \, dx$$

$$E(X) = \alpha \cdot \theta$$

$$\text{Var}(X) = \alpha \cdot \theta^2$$

19

# Example of stochastic models

12. Normal (Gaussian) distribution

$$X \sim N(\mu, \tau^2), \quad \theta = (\mu, \tau^2) \in \mathbb{R} \times (0, \infty)$$

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

$$E(x) = \mu$$

$$Var(x) = \tau^2$$

# Example of stochastic models

Quantitative continuous random vectors

13. Normal distribution $N(d; 0, I)$

- is the product of $d$ normal distributions $N(0,1)$

$$X = (x_1, \ldots, x_d)' \sim N(d; 0, I)$$

$$f(x) = \frac{1}{(2\pi)^{d/2}} \cdot \exp\left\{ -\frac{1}{2} x' x \right\}, \quad x = (x_1, \ldots, x_d)' \in \mathbb{R}^d$$

# Example of stochastic models

Quantitative continuous random vectors

13. Normal distribution $N(d; O, I)$
   - is the product of $d$ normal distributions $N(0,1)$

$$X = (x_1, \ldots, x_d)' \sim N(d; O, I)$$

$$f(x) = \frac{1}{(2\pi)^{d/2}} \cdot \exp\left\{-\frac{1}{2} x' x\right\}, \quad x = (x_1, \ldots, x_d)' \in \mathbb{R}^d$$

14. Normal distribution $N(d; \mu, \Sigma)$

$$X = (x_1, \ldots, x_d)' \sim N(d; \mu, \Sigma)$$

$$f(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} (\det \Sigma)^{\frac{1}{2}}} \cdot \exp\left\{-\frac{1}{2} (x-\mu)' \Sigma^{-1} (x-\mu)\right\},$$

$$x = (x_1, \ldots, x_d)' \in \mathbb{R}^d$$

$$E(X) = \mu$$

$$Cov(x, x') = E\left((x - E(x))(x - E(x))'\right) = \Sigma$$

$\Sigma$ is symmetric and positive definite matrix

# To do:

Write algorithms to generate the following random variables/vectors:

1) $X : \begin{pmatrix} a_1 & & a_m \\ p_1 & & p_m \end{pmatrix} \qquad \sum_{i=1}^{m} p_i = 1 \ , \ p_i \geq 0$

2) $X \sim Ge(\theta) \qquad \theta \in (0,1)$

3) $X = (X_1, \ldots, X_d)' \sim M(n \, ; p_1, \ldots, p_d) \qquad \sum_{i=1}^{d} p_i = 1 \ , \ p_i \geq 0$