

Warm-up sessions

Copernicus Master on Digital Earth

Probabilities and Statistics for data science

Prof. Nicolas Courty

ncourty@irisa.fr

Definitions

Definition of probability

Event space Ω All possible events together from a given experience.

Definition of $P(A)$ Let A be a set of events included in Ω ,

$$P(A) = \lim_{n \rightarrow \infty} \frac{n(A)}{n}$$

with

- n the number of experiments performed,
- $n(A)$ the number of experiments where A was performed.

Example, 6-sided dice

- $\Omega = \text{faces: } 1, 2, 3, 4, 5, 6$
- If dice not pipped, then $P(k) = 1/6, \quad \forall k \in 1, \dots, 6.$

Axioms of probabilities

- First axiom If $A \in \Omega$ then

$$0 \leq P(A) \leq 1$$

- Second axiom

$$P(\Omega) = 1 \quad P(\emptyset) = 0$$

with \emptyset the empty set

- Union and intersection If $A \in \Omega$, $B \in \Omega$, then

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

- If $A \cap B = \emptyset$ then

$$P(A \cup B) = P(A) + P(B)$$

Random Variable

Definition

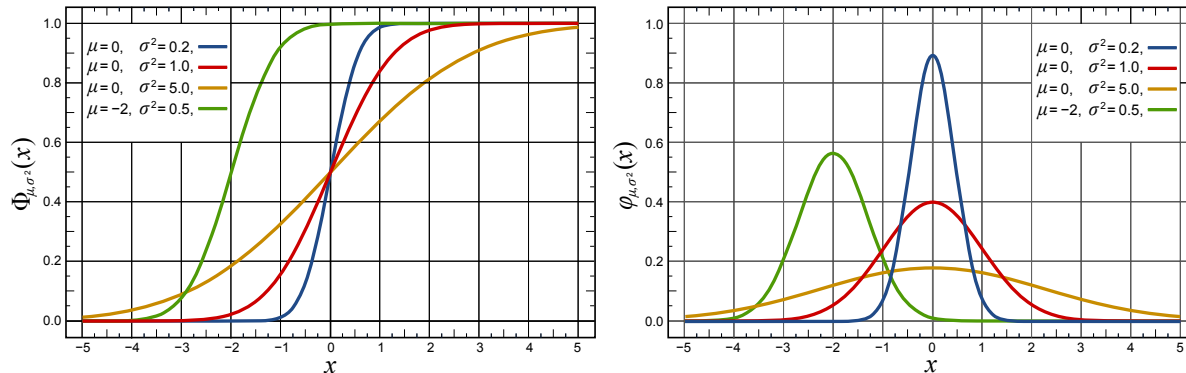
A random variable is a (real) X_ω number whose value is determined by the ω result of a randomized experiment.

Example: 6-sided dice

- The random event ω is the appearance of a face.
- An integer 1 to 6 is associated to each face.

Distribution function

Distribution function and derivative



- **Distribution function** (a.k.a. Cumulative distribution function) The F_X distribution function of a random variable (r.v.) X is defined as the probability that X is less than or equal to x ,

$$F_X(x) = P(X \leq x)$$

- **Probability density function.** It is defined as the derivative of the distribution function,

$$p(x) = \frac{dF(x)}{dx}$$

Properties

- Properties of the distribution function

$$F_X(-\infty) = 0, \quad F_X(\infty) = 1$$

$$0 \leq F_X(x) \leq 1$$

$$P(x_1 \leq x \leq x_2) = F_X(x_2) - F_X(x_1)$$

- Properties of the probability density

$$p(x) \geq 0 \quad \int_{-\infty}^{+\infty} p(x)dx = 1$$

$$P(x \leq x_1) = F_X(x_1) = \int_{-\infty}^{x_1} p(x)dx$$

$$P(x_1 \leq x \leq x_2) = \int_{x_1}^{x_2} p(x)dx$$

Moments of a random variable

Definition of a **moment**. The moment $g(x)$ of one random variable is given by its expectation

$$E(g(x)) = \int_{-\infty}^{+\infty} g(x)p(x)dx$$

Whenever $g(x) = x^m$, we refer to this quantity as the moment of order **m**,

$$\text{Moment of order 1} \quad m_X = E(X) = \int_{-\infty}^{+\infty} xp(x)dx$$

$$\text{Moment of order 2} \quad m_X^{(2)} = E(X^2) = \int_{-\infty}^{+\infty} x^2p(x)dx$$

The 1st order moment is also often called average (or mean).

- **Property** linearity of expectation

$$E(X + Y) = E(X) + E(Y), \quad E(kX) = kE(X)$$

For **k** a constant.

Moments of a random variable

Definition of the **variance**. The variance is the expectation of the square of the deviations from mean value $m_X = E(X)$,

$$\sigma_X^2 = E((X - m_X)^2) = \int_{-\infty}^{+\infty} (x - m_X)^2 p(x) dx ,$$
$$\sigma_X^2 = E(X^2) - E(X)^2$$

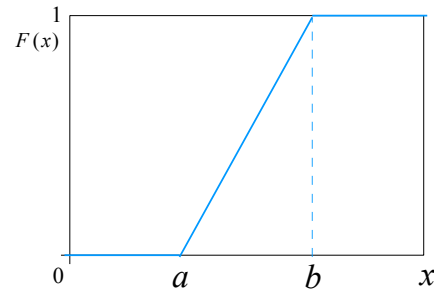
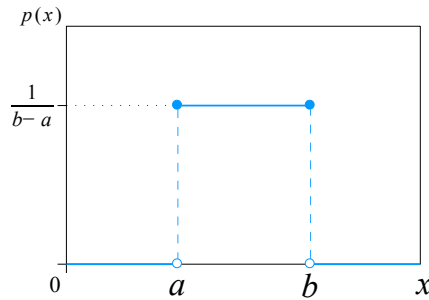
The notion of standard deviation σ is also often used,

$$\sigma_X = \sqrt{\sigma_X^2} .$$

- **Incomplete characterization**. Incomplete characterization of a random variable by its mean and variance.

Examples of laws

Uniform Law $\mathcal{U}(a, b)$



Probability density

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

Distribution function

$$F(x) = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & \text{if } x \in [a, b], \\ 1 & x' > b \end{cases}$$

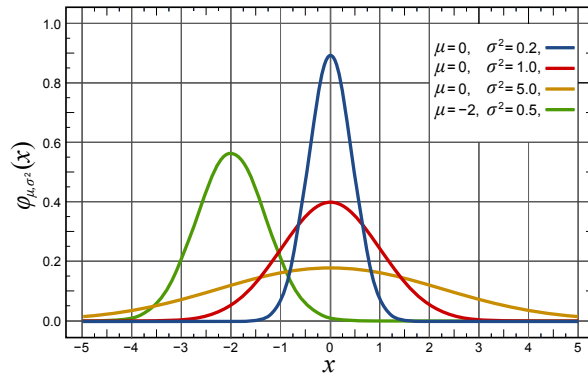
Expectation:

$$m_X = E(X) = \frac{b+a}{2}$$

Variance:

$$\text{Var}(X) = \frac{1}{12}(b-a)^2$$

Normal Law $\mathcal{N}(\mu, \sigma^2)$

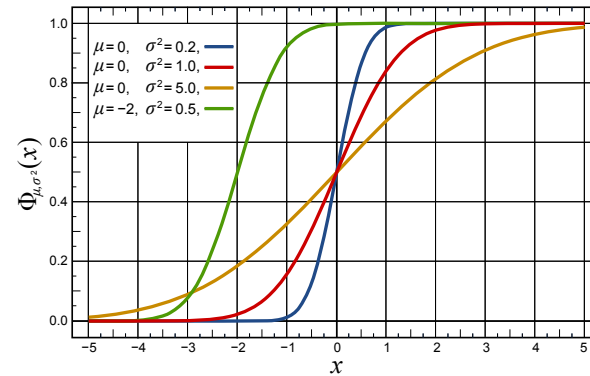


Probability density

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Distribution function

$$F(x) = \frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{x-\mu}{\sigma\sqrt{2}} \right) \right]$$



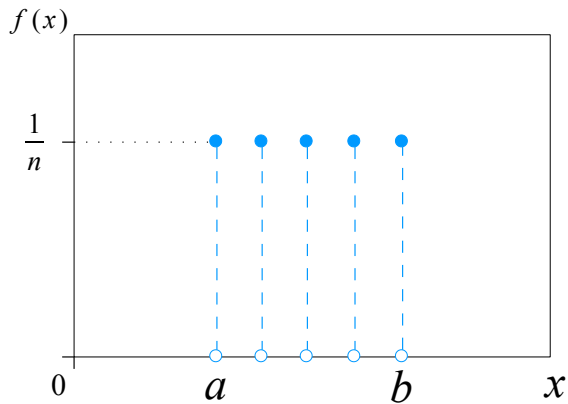
Expectation:

$$m_X = E(X) = \mu$$

Variance:

$$\operatorname{Var}(X) = \sigma^2$$

Empirical Law $\mathcal{U}(x_1, \dots, x_n)$

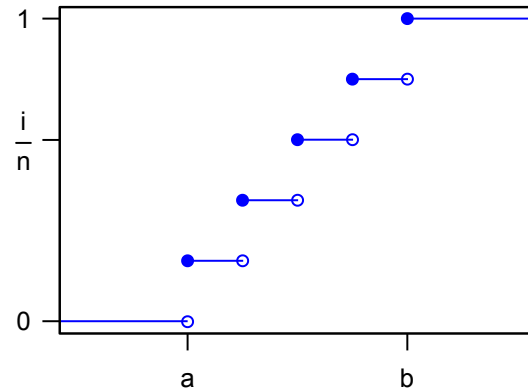


Probability density

$$p(x) = \frac{1}{n} \sum_{i=1}^n \delta(x - x_i)$$

Distribution function

$$F(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x \geq x_i}$$



Expectation:

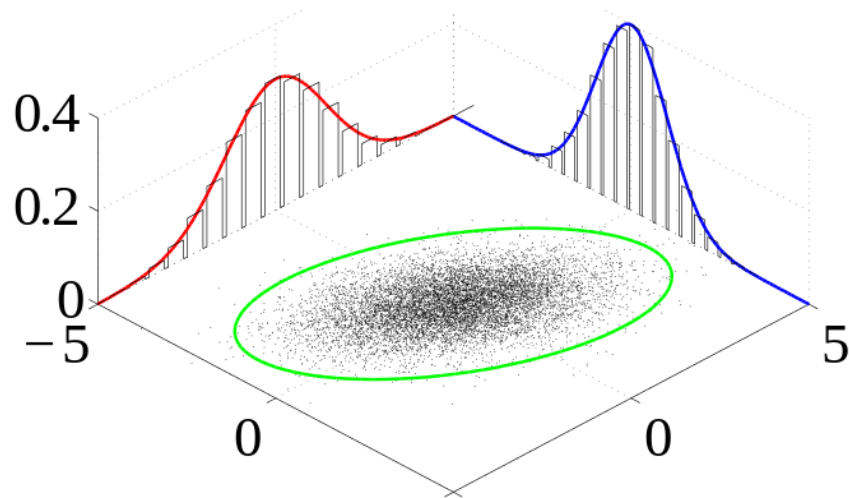
$$m_X = E(X) = \frac{1}{n} \sum_{i=1}^n x_i$$

Variance:

$$\text{Var}(X) = \frac{1}{n} \sum_{i=1}^n (x_i - m_X)^2$$

System of random variables

- We consider the case where several random variables are available simultaneously.
- We need to model jointly those variables



- When those variables X_1, X_2, \dots, X_d are given, we can model them by a vector of random variables $\mathbf{X} \in \mathbb{R}^d$

Joint probability density

Cumulative distribution function. Let

X and Y be two r.v. then,

$$F(x, y) = P(X \leq x, Y \leq y)$$

Joint probability density.

Let X and Y be two r.v. then,

$$p(x, y) = \frac{\partial^2 F(x, y)}{\partial x \partial y}$$

$$p(x) = \int p(x, y) dy$$

$$p(y) = \int p(x, y) dx$$

$p(x)$ and $p(y)$ are called marginal laws.

Properties

$$0 \leq F(x, y) \leq 1$$

$$F(-\infty, -\infty) = 0$$

$$F(\infty, \infty) = 1$$

Properties

$$p(x, y) \geq 0$$

$$\int p(x, y) dx dy = 1$$

$$p(A, B) = P(x \in A, y \in B)$$

$$= \int_A \int_B p(x, y) dx dy$$

Joint probability density

Conditional probability

- Joint law $p(\mathbf{x}, \mathbf{y})$.
- Probability of one of the variables knowing the value of the second.
- Notation: $p(\mathbf{x}|\mathbf{y})$.

Bayes' theorem

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{y})}$$

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})}$$

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$$



Covariance and correlation

- Moments of a joint law

$$E(g(x, y)) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g(x, y)p(x, y)dxdy$$

- Correlation

$$R_{XY} = E(XY) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} xyp(x, y)dxdy$$

- Covariance

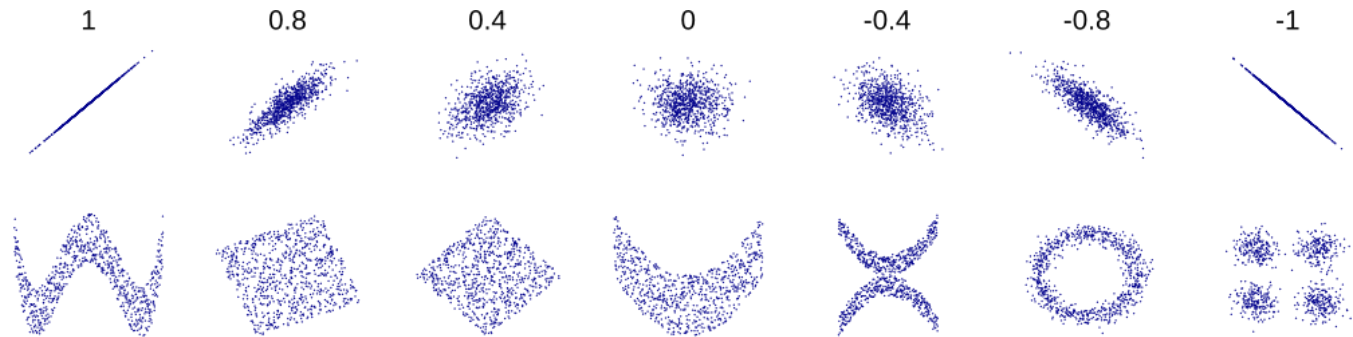
$$C_{XY} = \sigma_{XY} = E((X - m_X)(Y - m_Y))$$

$$C_{XY} = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} (x - m_X)(y - m_Y)p(x, y)dxdy$$

- Correlation coefficient

$$r_{XY} = \frac{C_{XY}}{\sigma_X \sigma_Y}$$

Independence and correlation



Covariance and Correlation

$$R_{XY} = E(XY) = C_{XY} + m_X m_Y$$

Independence

- Two r.v. X and Y are independent if

$$p(x, y) = p(x)p(y)$$

.

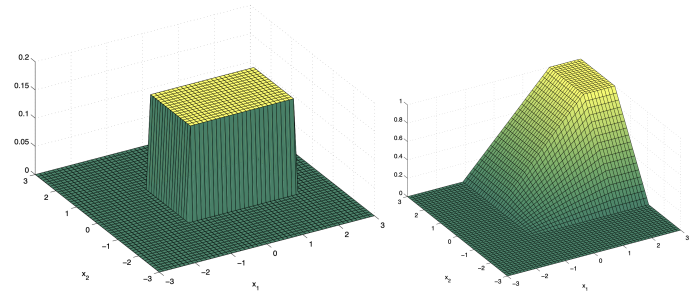
- If the variables are independent then

$$R_{XY} = m_X m_Y \quad \text{and} \quad C_{XY} = 0.$$

Examples of multivariate distribution function

Multivariate Uniform Law

- $X \sim U(a_x, b_x)$ and $Y \sim U(a_y, b_y)$
- $\mathbf{X} = [X, Y]^\top$



Probability density

$$p(x, y) = \begin{cases} \frac{1}{S} & \text{if } x \in [a_x, b_x] \\ & \text{and } y \in [a_y, b_y] \\ 0 & \text{else} \end{cases}$$

$$\text{Surface } S = (b_x - a_x)(b_y - a_y)$$

Expectation:

$$\mathbf{m}_X = E(\mathbf{X}) = \begin{bmatrix} \frac{b_x + a_x}{2} \\ \frac{b_y + a_y}{2} \end{bmatrix}$$

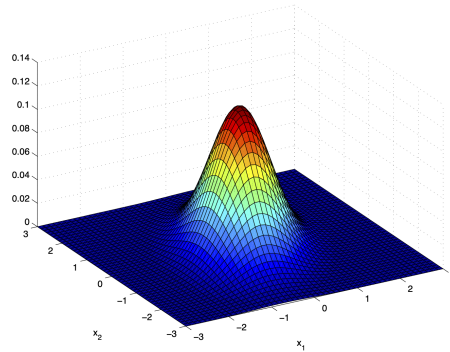
Variance:

$$\begin{aligned} \text{Cov}(\mathbf{X}) &= E((\mathbf{X} - \mathbf{m}_X)(\mathbf{X} - \mathbf{m}_X)) \\ &= \begin{bmatrix} \text{Var}(X) & 0 \\ 0 & \text{Var}(Y) \end{bmatrix} \end{aligned}$$

Examples of multivariate distribution function

Multivariate Gaussian Law

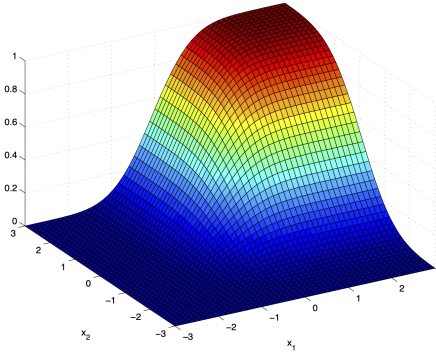
- $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$



Probability density

$$p(\mathbf{x}, \mathbf{y}) = K e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

$$\text{Coefficient } K = \frac{1}{(2\pi)^{N/2} |\boldsymbol{\Sigma}|^{1/2}}$$



Expectation:

$$\mathbf{m}_X = E(\mathbf{X}) = \boldsymbol{\mu}$$

Covariance:

$$\begin{aligned} \text{Cov}(\mathbf{X}) &= E((\mathbf{X} - \mathbf{m}_X)(\mathbf{X} - \mathbf{m}_X)^T) \\ &= \boldsymbol{\Sigma} \end{aligned}$$

Practical session

Let's practice now:

- Numpy, scipy and Pandas for manipulating dataframes
- Explore basic statistics computation using python

The end.

