

# **Statistical Data Analysis**

---

Dr. Jana de Wiljes

2. November 2021

Universität Potsdam

## **Continuous Random Variables**

---

# Normal Distribution

A normal or Gaussian distributed random variable  $X : \Omega \rightarrow \mathbb{R}$  with parameters  $\mu \in \mathbb{R}$  and  $\sigma > 0$  has the following density

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\}$$

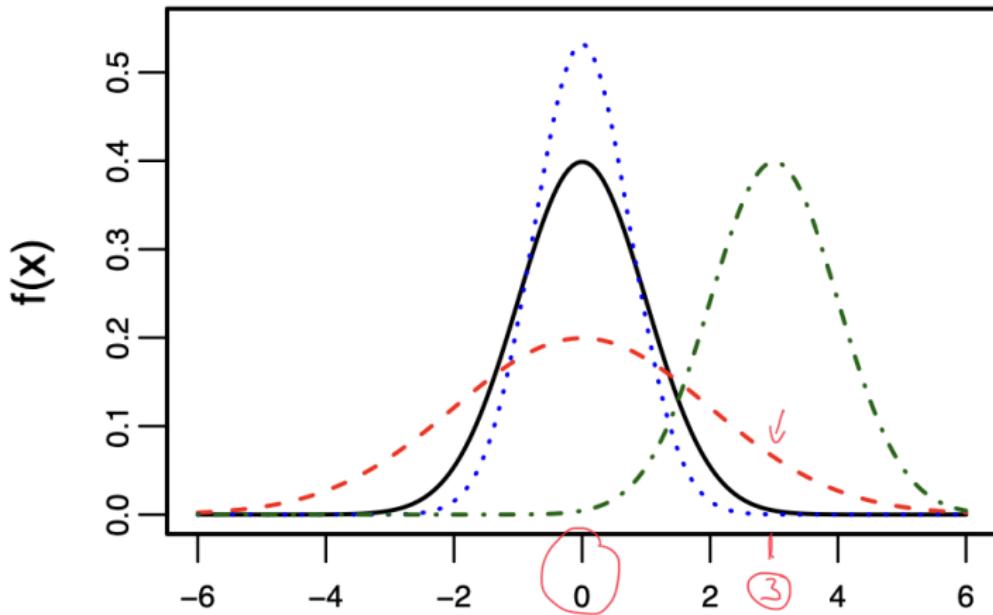
and expected value and variance

$$\mathbb{E}[X] = \mu$$

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

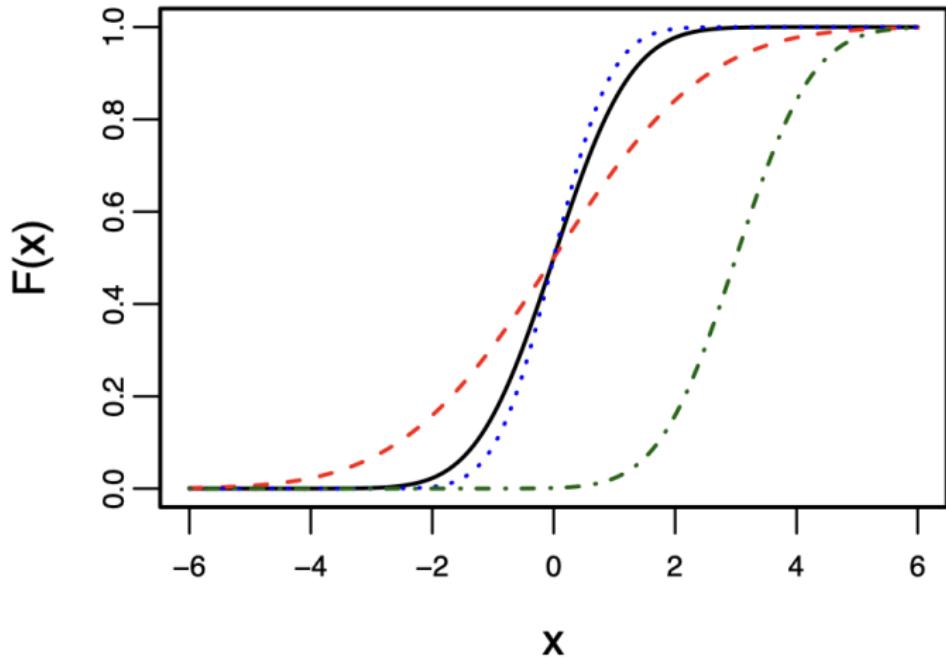
$$X \sim \mathcal{N}(\mu, \sigma^2)$$

# Normal Distribution



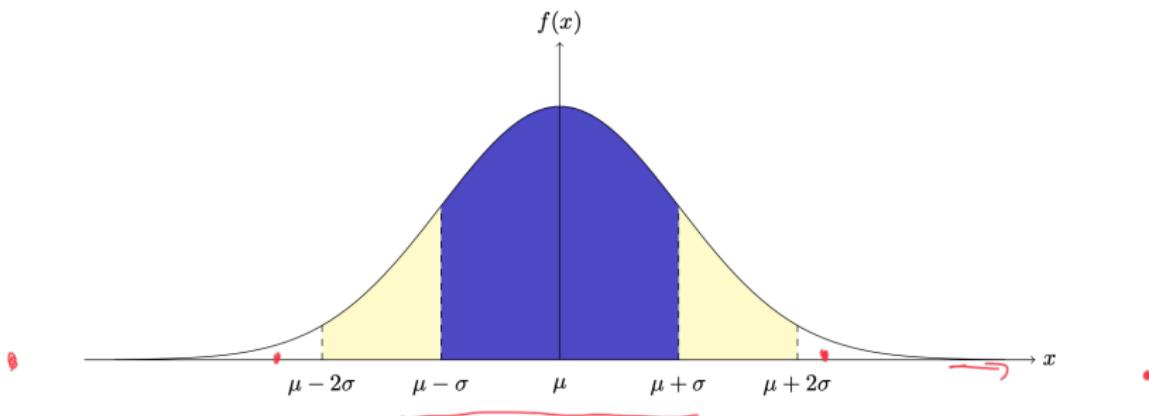
**Abbildung 1:**  $\mu = 0, \sigma = 1$  (black),  $\mu = 0, \sigma = 2$  (red),  $\mu = 0, \sigma = 0.75$  (blue) and  $\mu = 3, \sigma = 1$  (green)

# Normal Distribution



**Abbildung 2:**  $\mu = 0, \sigma = 1$  (black),  $\mu = 0, \sigma = 2$  (red),  $\mu = 0, \sigma = 0.75$  (blue) and  $\mu = 3, \sigma = 1$  (green)

# Quantile



68

**Abbildung 3:** ~~68~~<sup>68</sup>% of area under the curve (colored in blue) are in the  $[\mu - \sigma, \mu + \sigma]$  interval and 95% of the area under the curve are in the interval  $[\mu - 2\sigma, \mu + 2\sigma]$ .

99.7% in  $[\mu - 3\sigma, \mu + 3\sigma]$

## Standard normal distribution

A variable  $X : \Omega \rightarrow \mathbb{R}$  follows a standard normal distribution, i.e.,  $X \sim \mathcal{N}(0, 1)$  if the associated density has the following form

$$\phi(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\left(\frac{x^2}{2}\right)\right\}$$

with the associate cumulative distribution

$$\Phi(x) = \int_{-\infty}^x \phi(u) du \quad (1)$$

and quantile

$$z_\alpha = \Phi^{-1}(\alpha), \quad \alpha \in (0, 1) \quad (2)$$

Relationship between standard normal distribution and Normal distribution

$$F(x) = \Phi\left(\frac{x - \mu}{\sigma}\right) \quad (3)$$

# Exponential Distribution

A random variable  $X : \Omega \rightarrow \mathbb{R}$  follows the exponential distribution with parameters  $\lambda > 0$  has the following density and cdf

$$f(x) = \begin{cases} 0 & x < 0 \\ \lambda \exp(-\lambda x) & x \geq 0 \end{cases}$$

$$F(x) = \begin{cases} 0 & x < 0 \\ 1 - \exp(-\lambda x) & x \geq 0 \end{cases}$$

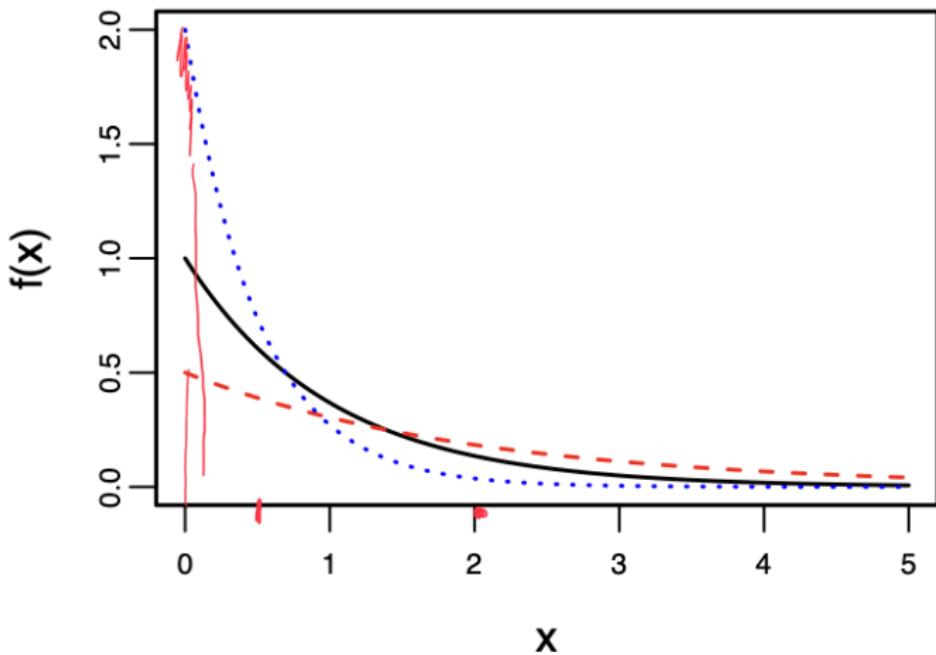
and expected value and variance

$$\mathbb{E}[X] = \frac{1}{\lambda}$$

$$\text{Var}(X) = \frac{1}{\lambda^2}$$

Notation:  $X \sim \text{Exp}(\lambda)$  (often used for waiting times and lifetimes)

# Exponential Distribution



**Abbildung 4:**  $\lambda = 1$  (black),  $\lambda = 2$  (blue) and  $\lambda = 1/2$  (red).

# Exponential Distribution

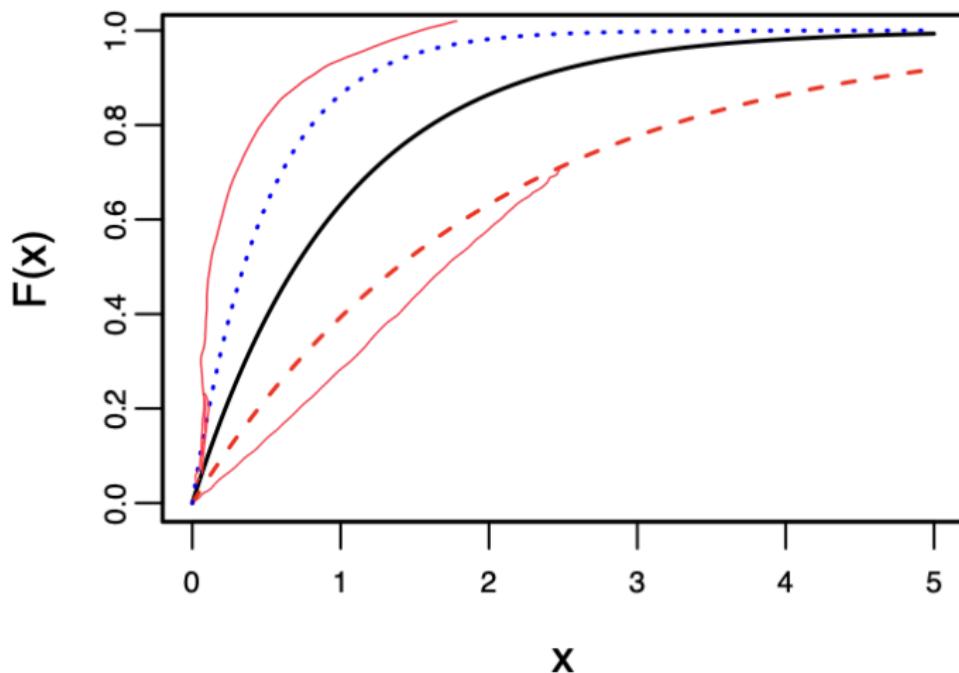


Abbildung 5:  $\lambda = 1$  (black),  $\lambda = 2$  (blue) and  $\lambda = 1/2$  (red).

## Example

**Setting:** The lifetime  $T$  of a computer chip is exponentially distributed, i.e.,  $T \sim \text{Exp}(\lambda)$  with expected lifetime of 15 weeks, i.e., parameter  $\lambda = \frac{1}{15}$

### Question:

- What is the probability that the computer chip is defect within the first 10 weeks?

$$P(T \leq 10) = F(10) = \underline{\underline{1 - e^{-\lambda \cdot 10}}} = 1 - e^{-\frac{10}{15}} = 0.487$$

- What is the probability that the computer chip will last at least 20 weeks?

$$P(T \geq 20) = 1 - P(T < 20) = \underline{\underline{1 - F(20)}} = e^{-\frac{20}{15}} = 0.264 \quad 10$$

## Transformation

Linear case:  $y$  linear i.e.,  $y(x) = bx + a$  for  $b > 0$ ,  $Y = g(x)$

$$\begin{aligned}F_Y(y) &= P(Y \leq y) = P(bx + a \leq y) = P\left(x \leq \frac{y-a}{b}\right) \\&= F_X\left(\frac{y-a}{b}\right)\end{aligned}$$

case  $b < 0$  we have

$$\begin{aligned}F_Y(y) &= P(Y \leq y) = P(bx \leq y-a) = P\left(X > \frac{y-a}{b}\right) \\&= 1 - F_X\left(\frac{y-a}{b}\right)\end{aligned}$$

# Transformation

On the level of the densities this boils down to:

for  $b \neq 0$  the density of  $Y$  with respect  $X$  is:

$$f_Y(y) = \left(\frac{1}{|b|}\right) \cdot f_X\left(\frac{y-a}{b}\right)$$

Example:  $X \sim N(\mu, \sigma^2)$      $y = a + bX \Rightarrow Y \sim N(a + b\mu, b^2\sigma^2)$   
with the density being

$$f_Y(y) = \frac{1}{\sqrt{2\pi}|\sigma|b} \exp\left\{-\frac{1}{2} \left(\frac{\frac{y-a}{b} - \mu}{\sigma}\right)^2\right\}$$

$f_X\left(\frac{y-a}{b}\right)$

# Transformation

---

# Transformation

**Reminder:** for arbitrary  $g$  the following holds:

$$\mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx \quad (4)$$

**Proposition:** Let  $g$  be a differentiable, strictly monotone function

and  $X$  a random variable. Then  $Y = g(X)$  has the following density

$$f_Y(y) = \left| \frac{1}{g'(g^{-1}(y))} \right| f_X(g^{-1}(y)), y \in E_Y \quad (5)$$

$E_Y$  is given by the value space of  $X$  via

$$E_Y = g(E_X) = \{g(x) : x \in E_X\} \quad (6)$$

## Example: Lognormal distribution

Example :  $X \sim N(\mu, \sigma^2)$  normal distribution

Then we call  $Y = e^X$  lognormal distribution

with

$$f_Y(y) = \begin{cases} 0 & y \leq 0 \\ \frac{1}{\sqrt{2\pi}\sigma y} \exp\left\{-\frac{1}{2}\left(\frac{\log(y)-\mu}{\sigma}\right)^2\right\} & y > 0 \end{cases}$$

## Jensen's inequality

**Proposition:** Let  $g$  be a convex function and  $X$  random variable

$$\mathbb{E}[g(X)] \geq g(\underline{\mathbb{E}[X]}) \quad (7)$$

**Example:**

For a lognormal-distribution  $Y$

$$\mathbb{E}[Y] = (e^{\mu + \sigma^2/2}) > g(\mu) = e^\mu$$

# Samples

**Definition:** Let  $(\Omega, \mathcal{F}, \mathbb{P})$  a probability space and  $\underline{X_1, \dots, X_n}$  be associated random variables. Realizations

$$x_1 := X_1(\omega), \dots, x_n := X_n(\omega) \quad (8)$$

are referred to as *samples* and  $n$  the sample size.

$$\begin{aligned} X_i : \Omega &\rightarrow E_{X_i} \\ \omega &\mapsto x_i \end{aligned}$$

# Estimator

**Definition:** A measurable function  $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is referred to as *sample function, estimator or statistic*.

Note: we will also consider the composition:

$$\varphi(X) : \Omega \rightarrow \mathbb{R}^m \quad (9)$$

$$\omega \mapsto \varphi(X_1(\omega), \dots, X_n(\omega)) \quad (10)$$

o

## Sample estimation

**Given:**  $(x_1, \dots, x_n) \in \mathbb{R}^n$  of independent and identical random variables  $X_1, \dots, X_n$  where

$$F(t) = \mathbb{P}[X_i \leq t], \quad t \in \mathbb{R} \quad (11)$$

but **unknown**

**Goal:** estimate  $\mathbb{E}[X_i]$  or  $\text{Var}[X_i]$



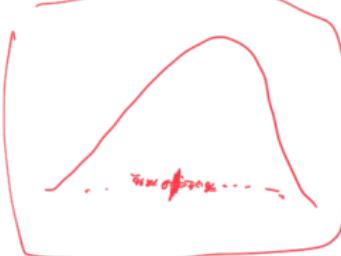
## Empirical mean

**Definition:** The empirical mean is defined by

$$\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i \quad (12)$$

Note: we will also use an analog notation for the random variables:

$$\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i \quad (13)$$



## Random variables

**Proposition:** Let  $X_1, \dots, X_n$  be independent and identical random variables with  $\mathbb{E}[X_i] = \mu$  and  $\text{Var}[X_i] = \sigma^2$ . Then

$$\mathbb{E}[\bar{X}_n] = \underline{\mu} \quad \text{and} \quad \text{Var}[\bar{X}_n] = \underline{\frac{\sigma^2}{n}} \quad (14)$$

# Proof

$$\bullet \mathbb{E}[\bar{X}_n] = \mathbb{E}\left[\frac{x_1 + \dots + x_n}{n}\right] \stackrel{\text{Def of } \bar{X}}{=} \frac{1}{n} \cdot \mathbb{E}[x_1 + \dots + x_n]$$

$\nwarrow \quad \nearrow$

$$= \frac{1}{n} \cdot n \cdot \underbrace{\mathbb{E}[x_i]}_{= \mu} = \mu \quad \square$$
$$\frac{1}{n} (\mathbb{E}[x_1] + \dots + \mathbb{E}[x_n])$$

$$\bullet \text{Var}(\bar{X}_n) = \text{Var}\left(\frac{x_1 + \dots + x_n}{n}\right) = \frac{1}{n^2} \cdot \text{Var}(x_1 + \dots + x_n)$$

Use this because  $x_i$  are independent

$$\frac{1}{n^2} \cdot n \cdot \text{Var}(x_i) = \frac{\sigma^2}{n}$$

□

$$\begin{aligned} \text{Var}(x+y) &= \text{Var}(x) + \text{Var}(y) \\ \text{if } x \text{ and } y \text{ are independent} \end{aligned}$$

# Proof

---

## Law of large numbers

---

**Proposition:** Let  $X_1, \dots, X_n$  be independent and identical random variables with  $\mathbb{E}[X_i] = \mu$ . Then

$$\bar{X}_n \rightarrow \mu \text{ for } n \rightarrow \infty \text{ (almost sure)} \quad (15)$$

## Empirical variance

---

**Definition:** The empirical variance is defined by

$$s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \quad (16)$$

Note: we will also use an analog notation for the random variables:

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \quad (17)$$

## Empirical variance

---

**Proposition:** Let  $X_1, \dots, X_n$  be independent and identical random variables. Then

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i^2 - n\bar{X}_n^2) \quad (18)$$

# Proof

---

# Proof

---

## Empirical variance

---

**Proposition:** Let  $X_1, \dots, X_n$  be independent and identical random variables with  $\mathbb{E}[X_i] = \mu$  and  $\text{Var}[X_i] = \sigma^2$ . Then

$$\mathbb{E}[S_n^2] = \sigma^2 \tag{19}$$

# Proof

---

# Proof

---