

# **Statistical Data Analysis**

---

Dr. Jana de Wiljes

9. November 2021

Universität Potsdam

## Empirical cdf of a sample set

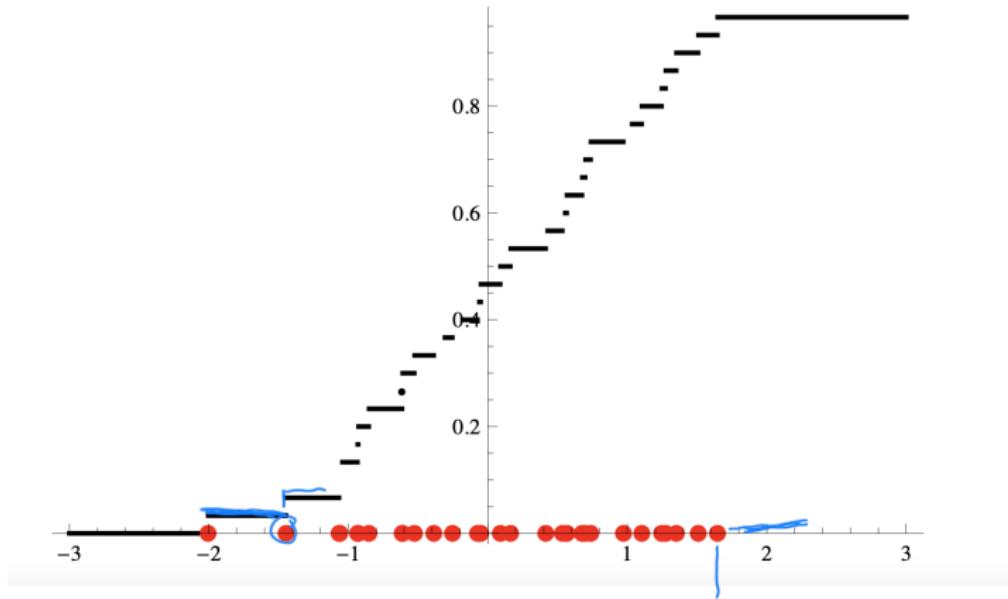
The empirical cdf of a sample set  $(x_1, \dots, x_n) \in \mathbb{R}^n$  is defined through

$$\widehat{F}_n(t) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i \leq t} = \frac{1}{n} \#\{i \in \{1, \dots, n\} : x_i \leq t\}, \quad t \in \mathbb{R} \quad (1)$$

Remark: the above defined empirical distribution can be approximated as follows via the order statistics  $x_{(1)}, \dots, x_{(n)}$

$$\widehat{F}_n(t) = \begin{cases} 0 & \text{for } t < x_{(1)} \\ 1/n & x_{(1)} \leq t < x_{(2)} \\ 2/n & x_{(2)} \leq t < x_{(3)} \\ \vdots & \vdots \\ n-1/n & \text{for } x_{(n-1)} \leq t < x_{(n)} \\ 1 & \text{for } x_{(n)} \leq t \end{cases}$$

## Empirical cdf



$$\hat{F}_n(t) \approx P(X \leq t)$$

$x_{(n)}$

## Empirical cdf

$X_1, \dots, X_n$  are i.i.d.

The empirical cdf of a sample set  $(x_1, \dots, x_n) \in \mathbb{R}^n$  is defined through

$$\widehat{F}_n(t) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{X_i \leq t} \quad (2)$$

Remark ②: The empirical distribution  $\widehat{F}_n$  has all the properties of cdf as following holds:

1)  $\lim_{t \rightarrow -\infty} \widehat{F}_n(t) = 0 \quad \lim_{t \rightarrow \infty} \widehat{F}_n(t) = 1$

2)  $\widehat{F}_n$  is monotonically increasing

3)  $\widehat{F}_n$  is right continuous

# Proposition

**Proposition:** Let  $(X_1, X_2, \dots)$  independent and identically distributed random variables with cdf  $F$ . Then

1 •

$$n\hat{F}_n(t) \sim \text{Bin}(n, F(t)).$$

This means

$$\mathbb{P}[\hat{F}_n(t) = \frac{k}{n}] = \binom{n}{k} F(t)^k (1 - F(t))^{n-k}, \quad k = 0, 1, \dots, n.$$

2 • The expect value and variance of  $\hat{F}_n(t)$  are given by

$$\mathbb{E}[\hat{F}_n(t)] = F(t), \quad \text{Var}[\hat{F}_n(t)] = \frac{F(t)(1 - F(t))}{n}$$

i.e.,  $\hat{F}_n(t)$  is an unbiased estimator of  $F(t)$ .

3 • For all  $t \in \mathbb{R}$  it holds that

$$\hat{F}_n(t) \xrightarrow{\text{a.s.}} F(t) \quad n \rightarrow \infty \text{ almost everywhere}$$

4 • For all  $t \in \mathbb{R}$  with  $F(t) \neq 0$  or 1 the following holds:

$$\sqrt{n} \frac{(\hat{F}_n(t) - F(t))}{\sqrt{F(t)(1 - F(t))}} \xrightarrow{\mathcal{N}(0, 1)} \text{for } n \rightarrow \infty \text{ (in distribution)}$$

# Proof

Remark: Part 4 can be interpreted as the distribution of the error

$\hat{F}_n(t) - F(t)$  is for large  $n$  approximately distributed according

$$N(0, \frac{F(t)(1-F(t))}{n})$$

Proof (7) Let us consider  $n$  experiments. For every  $i$ th experiment we check if  $X_i \leq t$ . In case  $X_i \leq t$  we say that the experiments a success, all the experiments are independent as the RV  $X_1, \dots, X_n$  are independent

$\Rightarrow$  the probability for success is  $P[X_i \leq t] = F(t)$

$\Rightarrow$  the number of successes in the  $n$  experiments is described by the random Variable

$$\hat{F}_n(t) = \sum_{i=1}^n \mathbf{1}_{X_i \leq t}$$

which is binomially distributed with parameters  $n$  and  $F(t)$

# Proof

Proof of (2) : we showed that  $n\hat{F}_n(t)$  distributed according to  $\text{Bin}(n, F(t))$

We know that the expected value and the variance of a random variable is that is binomial distributed (for  $n$  experiments and probability of success equal to  $F(t)$ )

$$\mathbb{E}[n\hat{F}_n(t)] = n \cdot F(t) \quad | \text{ taken out of } \mathbb{E}[\cdot] \text{ and divide by } n$$

$$\Rightarrow \mathbb{E}[\hat{F}_n(t)] = F(t)$$

The variance of a  $\text{Bin}(n, p)$  distributed random is

$$\text{Var}(X) = n \cdot p(1-p) \Rightarrow \text{Var}[n\hat{F}_n(t)] = n \cdot F(t) \cdot (1-F(t))$$

~ we can pull out  $n$  out of the variance and get a factor of  $n^2$  and then we divide by  $n^2$ . That yields

$$\text{Var}[\hat{F}_n(t)] = \frac{F(t) \cdot (1-F(t))}{n}$$

□

# Proof

Proof of (3) We introduce

$$Y_i = \begin{cases} 1 & X_i \leq t \\ 0 & \end{cases}$$

,  $Y_i$  are independent and identically distributed as  $X_i$  are.

$$\text{with } P[Y_i = 1] = P[X_i \leq t] = F(t)$$

$$P[Y_i = 0] = 1 - P[X_i \leq t] = 1 - F(t)$$

$\sim$  Bernoulli distributed

$$(E[Y_i] = F(t))$$

We can apply the law of large numbers to  $Y_1, Y_2, \dots$

$$\Rightarrow \hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \begin{cases} 1 & X_i \leq t \\ 0 & \end{cases} = \frac{1}{n} \sum_{i=1}^n Y_i \xrightarrow[n \rightarrow \infty]{\text{almost everywhere}} E[Y_i] = F(t)$$

$$\frac{1}{n} \sum X_i \rightarrow \mu$$

$E[X_i] = \mu$

□

# Proof

Proof of (4): Using the definition for  $Y_i$  (from proof of (3))

$$\mathbb{E}[Y_i] = F(t), \quad \text{Var}(Y_i) = F(t)(1-F(t))$$

(using the fact that  $Y_i$  are Bernoulli distributed)

And we make use of the central limit theorem for the series  $Y_1, Y_2, \dots$ .

$$\sqrt{n} \frac{\frac{1}{n} \sum_{i=1}^n (F(Y_i) - F(t))}{\sqrt{F(t)(1-F(t))}} = \sqrt{n} \frac{\sum_{i=1}^n Y_i - n\mathbb{E}[Y_i]}{\sqrt{\text{Var}(Y_i)}} = \frac{\sum_{i=1}^n Y_i - n\mathbb{E}[Y_i]}{\sqrt{n\text{Var}(Y_i)}}$$

convergence in probability  
 $\xrightarrow{n \rightarrow \infty}$

$$N(0, 1)$$

□

## Theoretical distribution

**Def:** Let  $X$  be a random. The theoretical distribution of  $X$  is a probability measure  $\mu$  on  $(\mathbb{R}, \mathcal{B})$  with

$$\mu(A) = \mathbb{P}[X \in A] \text{ for every Borel set } A \subset \mathbb{R} \quad (8)$$

**Note:** the relationship between the theoretical distribution  $\mu$  and the theoretical cdf  $F$  is;

$$F(t) = \mu((-\infty, t]), \quad t \in \mathbb{R} \quad (9)$$

## Empirical distribution

**Def:** The empirical distribution of a sample set  $(x_1, \dots, x_n) \in \mathbb{R}^n$  is defined through

$$\hat{\mu}_n(A) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i \in A} = \frac{1}{n} \#\{i \in \{1, \dots, n\} : x_i \in A\}, \quad (10)$$

for every Borel set  $A \subset \mathbb{R}$

$$\hat{F}(t) = \mu((-\infty, t])$$

$$[\circ, +]$$

## Dirac $\delta$ measure

**Def:** Let  $x \in \mathbb{R}$  be a real number. The dirac- $\delta$  measure  $\delta_x$  is a probability measure on  $(\mathbb{R}, \mathcal{B})$  with

$$\delta_x(A) = \begin{cases} 1, & \text{for } x \in A \\ 0, & \text{for } x \notin A \end{cases} \quad (11)$$

for all Borel set  $A \subset \mathbb{R}$



**Remark:** Then the empirical measure  $\hat{\mu}_n$  can be written as

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \quad (12)$$

and further note that

$$\hat{F}_n(t) = \hat{\mu}_n((-\infty, t]) \quad (13)$$

# Proposition

**Proposition:** Let  $(X_1, X_2, \dots)$  independent and identical distributed random variabel with distribution  $\mu$  and let  $A \subset \mathbb{R}$  a Borel set. Then

- 

$$n\hat{\mu}_n(A) \sim \text{Bin}(n, \mu(A)). \quad (14)$$

- The expect value and variance of  $\hat{\mu}_n(A)$  are given by

$$\mathbb{E}[\hat{\mu}_n(A)] = \mu(A), \quad \text{Var}[\hat{\mu}_n(A)] = \frac{\mu(A)(1 - \mu(A))}{n} \quad (15)$$

i.e.,  $\hat{\mu}_n(A)$  is an unbiased estimator of  $\mu(A)$ .

- Further it follows that  $\hat{\mu}_n$  is a consistent estimator, i.e.,

$$\hat{\mu}_n(A) \rightarrow \mu(A) \quad n \rightarrow \infty \text{ almost everywhere} \quad (16)$$

- For  $\mu(A) \neq 0$  or  $1$  the following holds:

$$\sqrt{n} \frac{\hat{\mu}_n(A) - \mu(A)}{\sqrt{\mu(A)(1 - \mu(A))}} \rightarrow \mathcal{N}(0, 1) \text{ for } n \rightarrow \infty \text{ (in distribution)} \quad (17)$$

# Plugin Estimator

**Setting:** Let  $(X_1, \dots, X_n)$  be independent and identically distributed random variables with the distribution  $\mu$ . Further we assume that a realisation  $(x_1, \dots, x_n) \in \mathbb{R}^n$  of the respective random variables

**Goal:** approximate  $\Psi(\mu)$  where  $\Psi : \mathcal{M} \rightarrow \mathbb{R}$

**Def:**  $\Psi(\hat{\mu}_n)$  is called the plugin estimator of  $\Psi(\mu)$ .

## Example

Example : (empirical moments). Let  $m_k := \Psi(\mu) = \int_{\mathbb{R}} x^k \cdot \mu(dx) = E[X^k]$  the  $k$ -th moment of  $\mu$ . The plugin estimator  $\hat{m}_k$  is given by

$$\hat{m}_k = \int_{\mathbb{R}} x^k \hat{\mu}_n(dx) = \frac{x_1^k + \dots + x_n^k}{n}$$

•  $\hat{m}_k$  is called  $k$ -th empirical moment of a sample set  $(x_1, \dots, x_n)$

Note that  $\boxed{\bar{X}_n = \hat{m}_1}$

---

Example: Let  $\Psi(\mu) = \text{Var } X_i = \int_{\mathbb{R}} x^2 \mu(dx) - \left( \int_{\mathbb{R}} x \mu(dx) \right)^2$

the variance associated with  $\mu$ . The plugin estimator is given by

$$\hat{\sigma}_{\text{plugin}}^2 = \int_{\mathbb{R}} x^2 \hat{\mu}_n(dx) - \left( \int_{\mathbb{R}} x \hat{\mu}_n(dx) \right)^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{X}_n^2 = \frac{n-1}{n} s_n^2$$

↑

where  $s_n^2$  is the empirical variance

$n-1$

## Example

---

## Kolmogorov-distance

**Def:** The Kolmogorov-distance between the empirical cdf  $\hat{F}_n(t)$  and the theoretical cdf  $F$  is defined as follows

$$D_n := \sup_{t \in \mathbb{R}} |\hat{F}_n(t) - F(t)| \quad (18)$$

# Theorem of Gliwenko-Cantelli

**Theorem:** For the Kolmogorov-distance  $D_n$  the following holds

$$D_n \xrightarrow{\text{a.s.}} 0 \text{ for } n \rightarrow \infty \text{ almost everywhere} \quad (19)$$

i.e.,

$$\mathbb{P}\left[\lim_{n \rightarrow \infty} D_n = 0\right] = 1 \quad (20)$$

Remark ⑥ Note that before we showed that  $\hat{F}_n(t) \rightarrow F(t)$  a.s.  
for every fixed  $t \in \mathbb{R}$

~ pointwise convergence

Now we show that it also converges uniformly

Remark ⑦ From the almost sure convergence ~~it follows~~ follows the convergence in distribution which means:

$$\lim_{n \rightarrow \infty} \mathbb{P}\left[\sup_{t \in \mathbb{R}} |\hat{F}_n(t) - F(t)| > \epsilon\right] = 0$$

# Proof

---

# Proof

---

# Proof

---

# Proof

---