

Statistical Data Analysis

Dr. Jana de Wiljes

9. November 2021

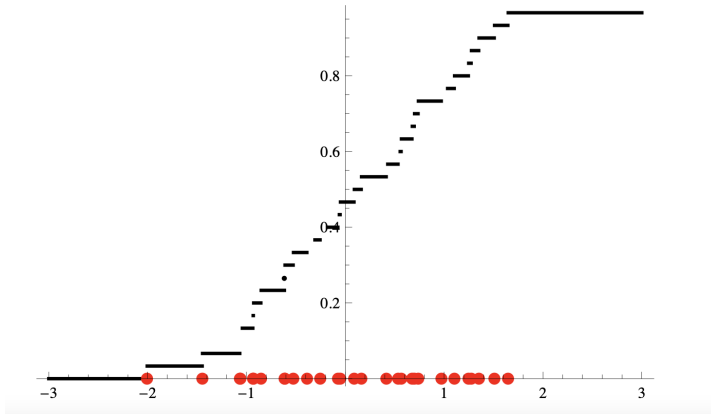
Universität Potsdam

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

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

Empirical cdf



Empirical cdf

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

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

Proposition

Proposition: Let (X_1, X_2, \dots) independent and identical distributed random variable with cdf F . Then

-

$$n\hat{F}_n(t) \sim \text{Bin}(n, F(t)). \quad (3)$$

This means

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

- The expected value and variance of $\hat{F}_n(t)$ are given by

$$\mathbb{E}[\hat{F}_n(t)] = F(t), \quad \text{Var}[\hat{F}_n(t)] = \frac{F(t)(1 - F(t))}{n} \quad (5)$$

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

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

$$\hat{F}_n(t) \rightarrow F(t) \quad n \rightarrow \infty \text{ almost everywhere} \quad (6)$$

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

$$\sqrt{n} \frac{\hat{F}_n(t) - F(t)}{\sqrt{F(t)(1 - F(t))}} \rightarrow \mathcal{N}(0, 1) \text{ for } n \rightarrow \infty \text{ (in distribution)} \quad (7)$$

Def: Let X be a random. The theoretical distribution of X is a probability measure μ 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 μ and the theoretical cdf F is;

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

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 $A \subset \mathbb{R}$

Dirac δ measure

Def: Let $x \in \mathbb{R}$ be a real number. The dirac- δ measure δ_x is a probability measure on $(\mathbb{R}, \mathcal{B})$ with 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 variable with distribution μ and let $A \subset \mathbb{R}$ a Borel set. Then

-

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

- The expected 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 identical distributed random variables with the distribution μ . 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

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

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

i.e.,

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

