

# ECDF and Plug-In Estimator

## Empirical Cumulative Distribution Function

Consider simple sample  $X_1, \dots, X_n \sim F(x)$ , where  $F(x)$  is unknown and we'd like to come up with an estimate  $\widehat{F}(x)$  of it. We would like that estimate to be **unbiased** and **consistent**.

Consider an estimate that is called the **empirical cumulative distribution function** (ECDF):

$$\widehat{F}(x) = \frac{1}{n} \sum_{k=1}^n \mathbb{I}nd\{X_n \leq x\}$$

where indicator function  $\mathbb{I}nd\{A\} = 1$  if event  $A$  is realized and 0 otherwise.

Is this estimate unbiased and consistent?

- $\mathbb{E} \left[ \widehat{F}(x) \right] = F(x)$ ?
- $\widehat{F}(x) \xrightarrow{P} F(x)$ ?

## ECDF properties

$$\widehat{F}(x) = \frac{1}{n} \sum_{k=1}^n \underbrace{\mathbb{I}nd\{X_n \leq x\}}_{\xi_k}$$

What is the distribution of  $\xi_k$ ?

By definition,  $\xi_k \sim Be(F(x))$  and, consequently,  $n\widehat{F}(x) \sim Bin(n, F(x))$ . What is the expected value and variance of binomial random variable?

- $\mathbb{E} \left[ \widehat{F}(x) \right] = \frac{1}{n} nF(x) = F(x)$  so ECDF is unbiased
- $\mathbb{V}ar \left( \widehat{F}(x) \right) = \frac{1}{n^2} nF(x)(1 - F(x)) \leq \frac{1}{4n} \rightarrow 0$  so ECDF is consistent

## ECDF convergence

We can estimate the speed of convergence from CLT:

$$\frac{\frac{1}{n} \sum_{k=1}^n \xi_k - \mathbb{E}[\xi_k]}{\sqrt{\mathbb{V}ar(\xi_k)}} \xrightarrow{d} \eta \sim \mathcal{N}(0, 1)$$

$$\sqrt{n} \left( \frac{1}{n} \sum_{k=1}^n \xi_k - \mathbb{E}[\xi_k] \right) \xrightarrow{d} \sqrt{\mathbb{V}ar(\xi_k)} \eta \sim \mathcal{N}(0, \mathbb{V}ar(\xi_k))$$

$$\sqrt{n} \left( \widehat{F}(x) - F(x) \right) \xrightarrow{d} \mathcal{N}(0, F(x)(1 - F(x)))$$

## ECDF uniform convergence

Glivenko-Cantelli theorem:

$$\sup_x \left| F(x) - \widehat{F}(x) \right| \xrightarrow{a.s.} 0$$

But how fast? **Kolmogorov's theorem**: If  $F(x)$  is continuous, then

$$D_n = \sqrt{n} \sup_x \left| F(x) - \widehat{F}(x) \right| \xrightarrow{d} \eta \sim K$$

$$\mathbb{P}(D_n \leq z) \rightarrow \sum_{k=-\infty}^{\infty} (-1)^k e^{-2k^2 z^2}$$

$$\mathbb{P}(D_n > z) \leq 2e^{-2nz^2}$$

## ECDF uniform coverage

If  $F(x)$  is continuous, then  $\sup_x \left| F(x) - \widehat{F}(x) \right|$  does not depend on  $F(\cdot)$ .

## Proof

$$\sup_x \left| \widehat{F}(x) - F(x) \right| = \sup_x \left| \frac{1}{n} \sum_{k=1}^n \mathbb{I}nd\{X_i \leq x\} - F(x) \right|$$

$$= \sup_x \left| \frac{1}{n} \sum_{k=1}^n \mathbb{I}nd\{F(X_i) \leq F(x)\} - \underbrace{F(x)}_{z \in [0,1]} \right|$$

$$= \sup_z \left| \frac{1}{n} \sum_{k=1}^n \mathbb{I}nd\{\underbrace{F(X_i)}_{\gamma} \leq z\} - z \right|$$

$$= \sup_z \left| \frac{1}{n} \sum_{k=1}^n \mathbb{I}nd\{U \leq z\} - z \right|$$

## Estimating functionals

A statistical functional  $T(F)$  is any functional of CDF  $F$ .

A **linear functional**  $T(F)$  can be written as:

$$T(F) = \int r(x) dF(x)$$

A **plug-in** estimator  $\widehat{T}$  of  $T(F)$  is

$$\widehat{T} = T(\widehat{F})$$

A plug-in estimator  $\widehat{T}$  of **linear**  $T(F)$  is

$$\widehat{T} = \int r(x) d\widehat{F}(x) = \frac{1}{n} \sum_{k=1}^n r(X_k)$$

## Estimating mean

Mean functional is:

$$\mu(F) = \int x dF(x)$$

So  $r(x) = x$ , and the plug-in estimator is

$$\widehat{\mu} = \frac{1}{n} \sum_{k=1}^n X_k = \overline{X}$$

## Estimating standard deviation

Standard deviation functional is

$$s(F) = \int (x - \mu)^2 dF(x)$$

So  $r(x) = (x - \mu)^2$  and the plug-in estimator is

$$\widehat{\sigma}_1^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \mu)^2$$

## Estimating standard deviation without the mean

$$\widehat{\sigma}_1^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \mu)^2$$

$$\widehat{\sigma}_2^2 = \frac{1}{n-1} \sum_{k=1}^n \left( X_k - \overline{X} \right)^2 = s^2$$

## Properties of normal distribution

Consider  $X_1, \dots, X_n \sim \mathcal{N}(m, \sigma^2)$ . Then,  $\widehat{\mu}$ ,  $\sigma_1^2$  and  $\sigma_2^2$  are unbiased consistent estimates.

## Proof 1

If  $X_1, \dots, X_n \sim \mathcal{N}(m, \sigma^2)$ , then  $\frac{1}{n} \sum_{k=1}^n X_k \sim ?$

$$\frac{1}{n} \sum_{k=1}^n X_k \sim \mathcal{N}\left(m, \frac{1}{n} \sigma^2\right)$$

So it is unbiased and consistent.

## Proof 2

$$\widehat{\sigma}_1^2 = \frac{1}{n} \sum_{k=1}^n (X_k - m)^2$$

$$\mathbb{E} \left[ \widehat{\sigma}_1^2 \right] = \mathbb{E} \left[ \frac{1}{n} \sum_{k=1}^n (X_k - m)^2 \right] = \frac{\sigma^2}{n} \mathbb{E} \left[ \sum_{k=1}^n \left( \frac{X_k - m}{\sigma} \right)^2 \right]$$

$$\eta = \sum_{k=1}^n \left( \frac{X_k - m}{\sigma} \right)^2 \sim ?$$

$$\eta = \sum_{k=1}^n \left( \frac{X_k - m}{\sigma} \right)^2 \sim \chi^2(n)$$

$$\mathbb{E}[\eta] = ?$$

$$\mathbb{V}ar(\eta) = ?$$

## Proof 2

$$\eta = \sum_{k=1}^n \left( \frac{X_k - m}{\sigma} \right)^2 \sim \chi^2(n)$$

$$\mathbb{E}[\eta] = n$$

$$\mathbb{V}ar(\eta) = 2n$$

$$\mathbb{E} \left[ \widehat{\sigma}_1^2 \right] = \frac{\sigma^2}{n} \mathbb{E}[\eta] = \frac{\sigma^2}{n} n$$

$$\mathbb{V}ar \left( \widehat{\sigma}_1^2 \right) = \frac{\sigma^4}{n^2} \mathbb{V}ar(\eta) = \frac{\sigma^4}{n^2} 2n$$