

ESTIMATING THE CDF

Let X_1, \dots, X_n be iid with continuous cdf F .

Defn: The Empirical cdf is the function

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \leq x), \quad x \in \mathbb{R}.$$

Result: For a fixed $x_0 \in \mathbb{R}$ we have

$$\sqrt{n}(\hat{F}_n(x_0) - F(x_0)) \xrightarrow{D} N(0, F(x_0)[1 - F(x_0)]).$$

Proof: Note that

$$\mathbb{E} \mathbb{1}(X_1 \leq x_0) = F(x_0)$$

$$\text{Var} \mathbb{1}(X_1 \leq x_0) = F(x_0)[1 - F(x_0)],$$

Then the result follows directly from the Central Limit Theorem.

Result: For F any continuous cdf, we have

$$(Wass, pg 14) (i) \quad \sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)| \rightarrow 0 \quad \text{u.p.} \quad (\text{Glivenko-Cantelli})$$

$$(Wass, pg 14) (ii) \quad P\left(\sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)| \leq \varepsilon\right) \geq 1 - 2e^{-2n\varepsilon^2} \quad (\text{DKW})$$

$$[\text{See Massart (1990)}] (iii) \quad \sqrt{n} \left(\sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)| \right) \xrightarrow{D} \sup_{t \in [0,1]} |B_0(t)|,$$

where $\{B_0(t) : t \in [0,1]\}$ is a Brownian bridge (defined below).

Defn: A Wiener process B is a random function in the space $C[0,1]$ of continuous functions on $[0,1]$ which satisfies
(Serf, pg 14)

(a) $B(0) = 0$ with probability 1.

(b) $B(t) \sim N(0, t)$, $t \in (0, 1]$

(c) For $0 \leq t_0 \leq t_1 \leq \dots \leq t_k \leq 1$, the increments $B(t_1) - B(t_0), \dots, B(t_k) - B(t_{k-1})$

are mutually independent.

Also called "Standard Brownian Motion."

Result: (Allows us to simulate Brownian motion) For each $n \geq 1$, let

$$B_n(t) = \frac{1}{\sqrt{n}} \sum_{i=1}^{\lfloor tn \rfloor} Z_i, \quad Z_1, \dots, Z_n \stackrel{\text{ind}}{\sim} N(0, 1).$$

Then B_n converges to B by a theorem called Donsker's Thm.

Defn: A Brownian bridge is the random function in $C[0,1]$ given by
(Athla, pg 375)

$$B_0(t) = B(t) - tB(1),$$

where B is a standard Brownian motion.

Remark: The "bridge" begins and ends at $0 = B_0(0) = B_0(1)$.

Exercise: Build confidence bands for F :

- with DFW result.
- with KS result (simulate to get needed quantile). 12

Solution: To use the DKW result, we write

$$P\left(\sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)| \leq \varepsilon\right) \geq 1 - 2e^{-2n\varepsilon^2} = 1 - \alpha$$

Now we have

$$1 - 2e^{-2n\varepsilon^2} = 1 - \alpha$$

\Leftrightarrow

$$-2n\varepsilon^2 = \log\left(\frac{\alpha}{2}\right)$$

\Leftrightarrow

$$\varepsilon = \sqrt{\frac{\log(2) - \log(\alpha)}{2n}}$$

So $(1 - \alpha)100\%$ confidence bands may be constructed as

$$\hat{F}_n(x) \pm \sqrt{\frac{\log(2) - \log(\alpha)}{2n}} \quad \text{for all } x \in \mathbb{R}.$$

From the KS result we could construct the interval

$$\hat{F}_n(x) \pm \frac{1}{\sqrt{n}} B_{0, \alpha/2} \quad \text{for } x \in \mathbb{R},$$

where $B_{0, \alpha/2}$ represents the upper $\alpha/2$ quantile of $\sup_{t \in [0, 1]} |B_0(t)|$.

Result: For a Brownian bridge B_0 , we have

$$\left[\begin{array}{c} \text{See} \\ \text{Massart (1990)} \end{array} \right] P\left(\sup_{t \in [0, 1]} |B_0(t)| \leq x\right) = 1 - 2 \sum_{i=1}^{\infty} (-1)^{i+1} \exp(-2i^2 x^2).$$

Exercise: Check accuracy of simulated quantiles of $\sup_{t \in [0, 1]} |B_0(t)|$.