

THE EMPIRICAL DISTRIBUTION

STAT 510

SETUP

$$X_1, X_2, \dots, X_n \sim F$$

- How TO ESTIMATE F ?
- How TO ESTIMATE PROPERTIES OF F ?

DEFN THE EMPIRICAL DISTRIBUTION FUNCTION, \hat{F}_n ,

IS THE CDF THAT PUTS MASS $\frac{1}{n}$ AT EACH POINT X_i :

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$$

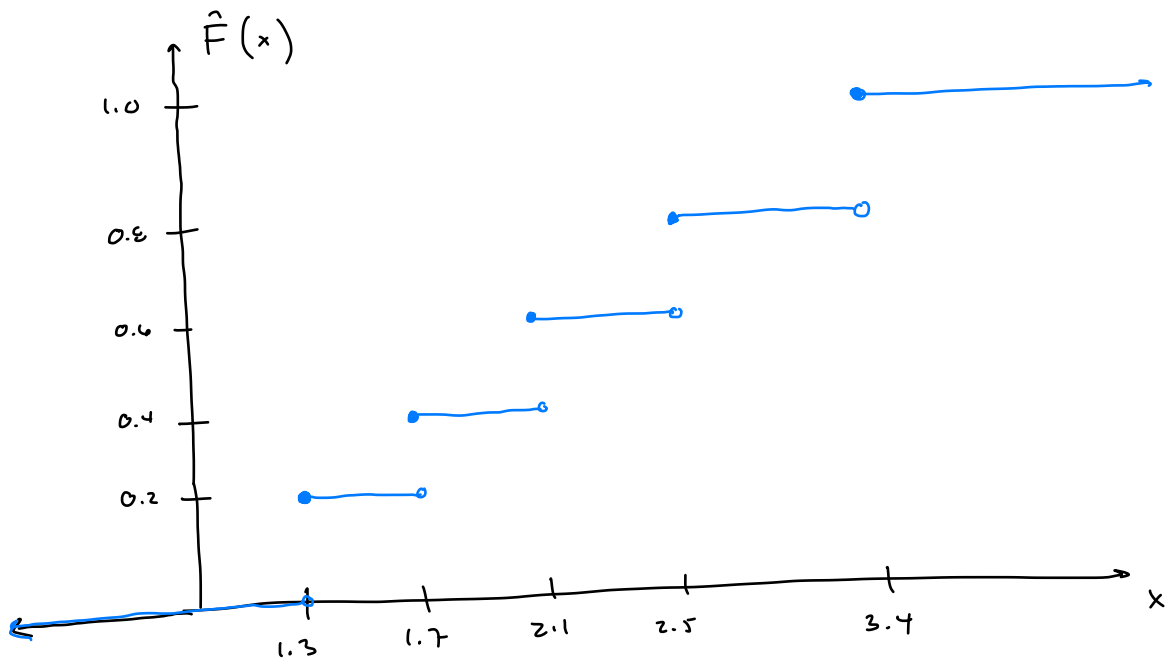
WHERE

$$I(X_i \leq x) = \begin{cases} 1 & \text{IF } X_i \leq x \\ 0 & \text{IF } X_i > x \end{cases}$$

EXAMPLE

GIVEN DATA 2.1, 1.3, 3.4, 2.5, 1.7
SORTED: 1.3, 1.7, 2.1, 2.5, 3.4

$$\hat{F}_n(x) = \begin{cases} 0 & x < 1.3 \\ 0.2 & 1.3 \leq x < 1.7 \\ 0.4 & 1.7 \leq x < 2.1 \\ 0.6 & 2.1 \leq x < 2.5 \\ 0.8 & 2.5 \leq x < 3.4 \\ 1 & x \geq 3.4 \end{cases}$$



Thm

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq x)$$

FOR ANY FIXED x

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

$$\mathbb{V}[\hat{F}_n(x)] = \frac{F(x)(1-F(x))}{n}$$

$$\text{MSE} = \frac{F(x)(1-F(x))}{n} \rightarrow 0$$

$$\hat{F}_n(x) \xrightarrow{P} F(x)$$

THM GLIVENKO - CANTELLI THEOREM

$$X_1, \dots, X_n \sim F$$

THEN

$$\sup_x |\hat{F}_n(x) - F(x)| \xrightarrow{P} 0$$

Thm DKW INEQUALITY

$$X_1, \dots, X_n \sim F$$

FOR ANY $\epsilon > 0$

$$P\left(\sup_x |\hat{F}_n(x) - F(x)| > \epsilon\right) \leq 2e^{-2n\epsilon^2}$$

NONPARAMETRIC $1-\alpha$ CONFIDENCE BAND FOR F

DEFINE $\epsilon_n = \sqrt{\frac{1}{2n} \log\left(\frac{2}{\alpha}\right)}$

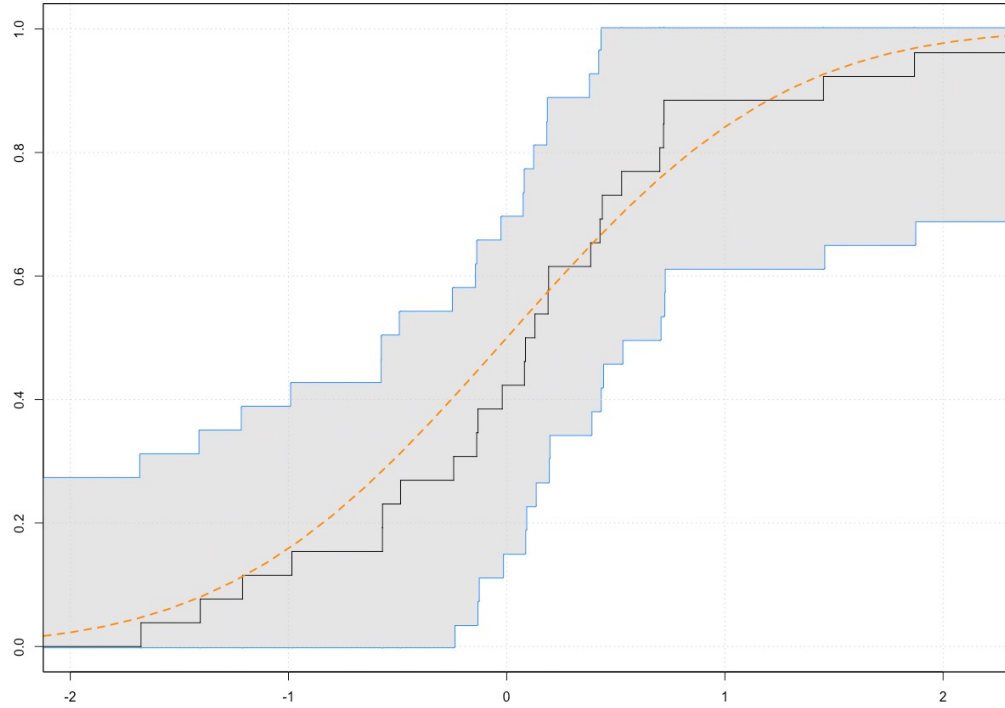
$$L(x) = \max\left\{\hat{F}_n(x) - \epsilon_n, 0\right\}$$

$$U(x) = \min\left\{\hat{F}_n(x) + \epsilon_n, 1\right\}$$

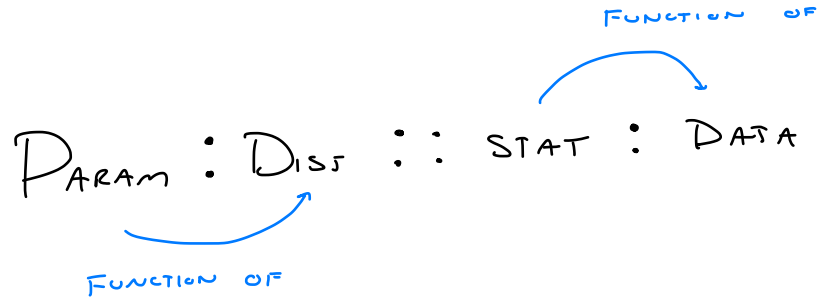
THEN DKW GIVES

$$P\left(L(x) \leq F(x) \leq U(x) \text{ FOR ALL } x\right) \geq 1 - \alpha.$$

$X_i \sim N(0,1) \quad n=25$



IDEA



STATISTICAL FUNCTIONALS

ANY FUNCTION $T(F)$ OF F .

• $\mu = \int x dF(x)$ MEAN

• $\sigma^2 = \int (x - \mu)^2 dF(x)$ VARIANCE

• $m = F^{-1}(1/2)$ MEDIAN

DEFN

PLUG-IN ESTIMATOR OF $\theta = T(F)$

$$\hat{\theta}_n = T(\hat{F}_n)$$

"PLUG-IN" \hat{F}_n FOR F !

DEFN

$$\text{IF } T(F) = \int r(x) dF(x) \quad \text{FOR SOME } r(x)$$

THEN T IS A LINEAR FUNCTIONAL.

$$T(aF + bG) = aT(F) + bT(G)$$

Talm

$$T(\hat{F}_n) = \int r(x) d\hat{F}_n(x) = \sum_{i=1}^n r(x_i) f(x_i) = \frac{1}{n} \sum_{i=1}^n r(x_i)$$

Annotations:

- An arrow points from \hat{F}_n to the text "DISCRETE".
- An arrow points from $\frac{1}{n}$ to the text "EACH POINT HAS PROB $1/n$ ".
- An arrow points from $T(\hat{F}_n)$ to the text "DEPN".

$$\text{IF } T(\bar{F}_n) \approx N(T(F), \hat{s}e^2)$$

THEN A 95% "NORMAL" CI FOR $T(F)$ IS APPROX

$$T(\bar{F}_n) \pm 2 \hat{s}e$$

EXAMPLE

MEAN

$$\text{LET } \mu = T(F) = \int x dF(x)$$

$$\hat{\mu} = T(\hat{F}_n) = \int x d\hat{F}_n(x) = \frac{1}{n} \sum x_i = \bar{x} \quad \checkmark$$

WE OFTEN KNOW $V[\bar{x}] = \sigma^2/n$

EXAMPLE

VARIANCE

$$\text{LET } \sigma^2 = T(F) = V[X] = \int (x-\mu)^2 dF(x)$$

$$= \int x^2 dF(x) - \left(\int x dF(x) \right)^2$$

$$\hat{\sigma}^2 = \int x^2 d\hat{F}_n(x) - \left(\int x d\hat{F}_n(x) \right)^2$$

$$= \frac{1}{n} \sum x_i^2 - \left(\frac{1}{n} \sum x_i \right)^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

$$V[\hat{\sigma}^2] = ???$$

$$\neq \frac{1}{n-1} \sum (x_i - \bar{x})^2 = S_n^2$$