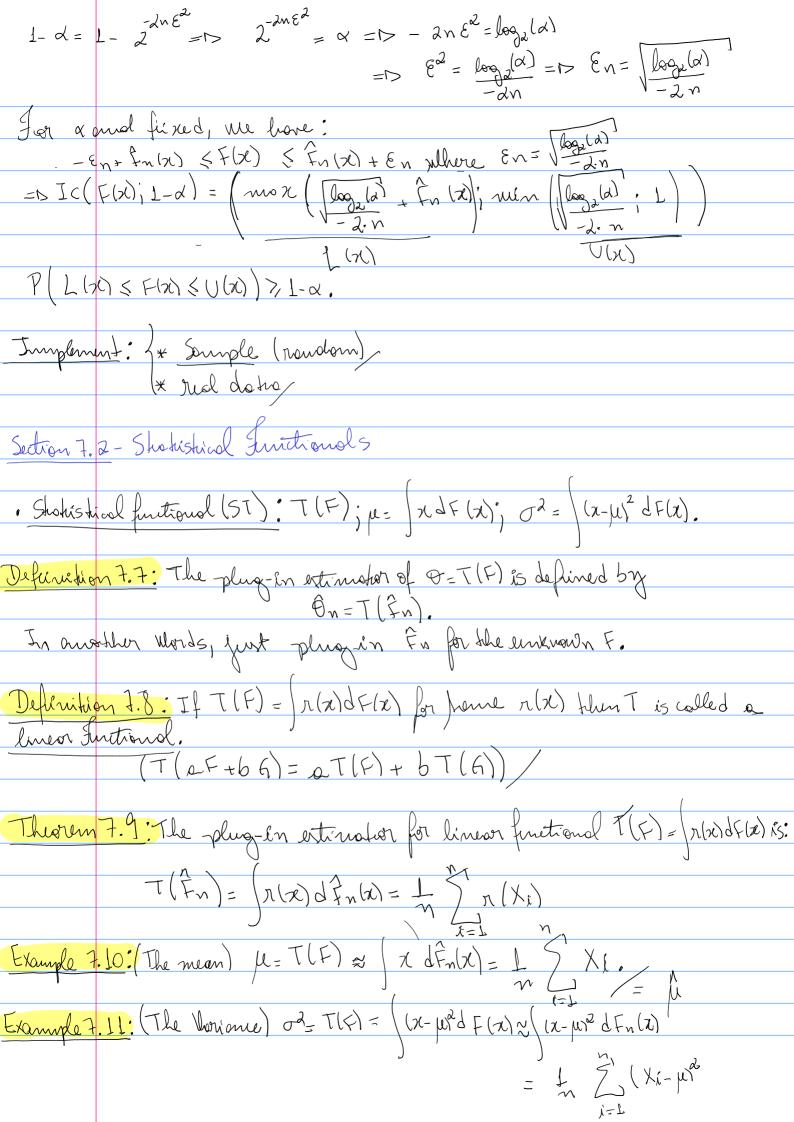
Chapter 7	- Esternaturg CDF and Estatistical Internal
Sechion 7.	L- The Emphrical Distribution hundren XI,, Xn ind F
Defunchio	at each date point $\chi_i$ . Formally, $\hat{f}_n(x) = \sum_{k=1}^{n} I(\chi_k \leq x)$ ,
mars In	at each date point Xi Formally,
	$f_n(x) = \frac{\sum_{k=1}^n \frac{1}{2}(x_k \leq x)}{x_k}$
1.10 T/	X (-a)   1 an y >
Where I	$(X_i \leq \mathcal{R}) = \begin{cases} 1, & \text{if } X_i > \mathcal{R}, \\ 0, & \text{if } X_i < \mathcal{R}. \end{cases}$
Theorem	13° Alanu Rived amlie Ala.
NOVO V	· In (21) Ph F(X)
	· E (Fh(x)) = F(x)
	$Vor(F_n(x)) = F(x)(1-F(x))$
	$\gamma$
	MSE = F(x)(1-F(x)) -> 0
T 1.	1 . l /
- myll r	nentil/
The ending on 7	49 The Gliversko-Contulli Theorem) Lot X. X. iid F. Theorem
10000117	42 The Glivenko-Contelli Theorem) Let XI, -, Xn ind F. Then  My I Fn(x) - F(x)   a.s. D. O.
Theorem	7.5: (The Drove to ky- keefer - Walfoulitz (D KW) I regurding)  n tid F, then for any 6>0, we have
fet Xn. X	n tid F, then for any 6>0, we have
	$\mathcal{D}(1)$ $\mathcal{C}(1)$ $\mathcal{C}(2)$ $\mathcal{C}(3)$ $\mathcal{C}(3)$ $\mathcal{C}(3)$ $\mathcal{C}(3)$ $\mathcal{C}(3)$
	$P\left(\left \sup_{\mathcal{H}}\left \widehat{F}_{n}\left \chi\right -F\left(\chi\right)\right >\varepsilon\right)\leqslant2^{-2n\varepsilon^{2}}$
P ( Jug	Fn (x) - F(x) > E) + P ( Imp   Fn(x) - F(x)   S E) = 1
70	,
1 M	p (Fn (x) - F(x)) ( E) = 1 - P ( Mp (Fn (x) - F(x)) > e)
	>1-2nEd
	Fn(x)-F(x)(SE) C[17.1x)- E(x)(SE)
1-02	
1-2°	$  F_{n}(x) - F(x)   \le \varepsilon   C  ^{2}   F_{n}(x) - F(x)   \le \varepsilon   F_{n}(x) - F(x) $



Example 7.12: (The Shewness) - $K = \frac{\int (x - \mu)^3 dF(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF(x)\right)^{\frac{3}{2}}} \sim \frac{\int (x - \mu)^3 dF_n(x)}{\left[\left((x - \mu)^2 dF$ Example 7.14: (Quantiles).  $T(F) = F^{-}(p) = D T(F) \approx \widehat{F}_{n}^{-}(p) / F_{n}(p) = enf$   $f_{n}(x) > p$ .