

Konf. int. / fornto.
 X_1, X_2, \dots, X_n uavh $N(\mu, \sigma^2)$?
 $T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim T_{n-1}$
 $\rightarrow \bar{X} \pm t_{\alpha/2} \frac{S}{\sqrt{n}} =$ eksakt $(1-\alpha)100\%$ KI for μ .
 $\frac{S^2}{\sigma^2} (n-1) \sim \chi_{n-1}^2$
 $\rightarrow (\frac{S^2}{s_{\alpha/2}^2} (n-1), \frac{S^2}{s_{1-\alpha/2}^2} (n-1))$ er et eksakt $(1-\alpha)100\%$ KI for σ^2
 $P(\chi_{n-1}^2 > \chi_{\alpha}^2) = \alpha$
 T-intervallet er robust for normalitetsantagelsen når n er stor
 Kji-int. krever normalfordeling også ved stor n .
 Men når n er liten og/eller normalfordeling ikke holder er det ikke sikkert at
 Dekningsgrad = Coverage = $P(\theta \in (\theta_L, \theta_U)) = 1-\alpha$

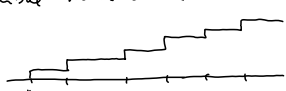
sep 12-12:13

Hvis antagelsen riktigen kan vi istedet beregne KI ved bootstrapping.
Bootstrap Data $X_1, \dots, X_n \sim f(x; \theta)$
 (0) Estimat $\hat{\theta}$ fra X_1, \dots, X_n
 (1) Generer pseudo data X_1^*, \dots, X_n^*
 (2) Beregn θ^* fra pseudo dataene
 (3) Gjenta (1) - (2) B ganger, $\theta_1^*, \dots, \theta_B^*$
 Hvis fordelingen til θ_i^* -ene tiln. normal får et tiln. KI ved
 $\hat{\theta} \pm t_{\alpha/2} S_{\theta^*}$ Boot-t interval
 der $S_{\theta^*}^2 = \frac{1}{B} \sum_{b=1}^B (\theta_b^* - \bar{\theta}^*)^2$

sep 12-12:44

Alternativt kan vi benytte peranti (intervallet)
 Sorter θ_i^* -ene: $\theta_{(1)}^* < \theta_{(2)}^* < \dots < \theta_{(n)}^*$
 La f.eks. $b_L = [0.025B]$ (og $[x]$ = heltallsverdi) \leftarrow KI
 $b_U = [0.975B]$
 Da er $(\theta_{(b_L)}^*, \theta_{(b_U)}^*)$ et 95% peranti-interval for θ .
Parametrisert bootstrap:
 Pseudo data X_1^*, \dots, X_n^* trekkes fra $f(x; \hat{\theta})$
Ikke-parametrisert bootstrap:
 Generer X_1^*, \dots, X_n^* ved å trekke med tilbakelegging for X_1, \dots, X_n

sep 12-12:51

Ikke-parametrisert bootstrap
 er å trekke fra fordeling $P(X^* = x_i) = \frac{1}{n}$ observerte x_1, \dots, x_n (hvis alle ulike)
 Generelt: Trekk fra empiriske kumulative $\hat{F}_n(x) = \frac{1}{n} \# \{x_i \leq x\}$

 Hvis at $\hat{F}_n(x) \rightarrow F(x) = P(X \leq x)$ $n \rightarrow \infty$
 $E[\hat{F}_n(x)] = F(x)$
 $Var[\hat{F}_n(x)] = \frac{1}{n} F(x)(1-F(x))$
 Med boot-t interv. trenger vi bare est. av $\hat{\theta}$, kan holde med $B = 100$
 Men for peranti intervallet bør $B = 1000$ (est. store)

sep 12-12:58

Teoretisk median M er def. ved
 $P(X \leq M) = \frac{1}{2} = P(X \geq M)$
 (når X er kont. fordelt)
 Median: utvalgt $m = \begin{cases} X_{m+1} & \text{når } n = 2m+1 \\ \frac{X_m + X_{m+1}}{2} & \text{når } n = 2m \end{cases}$
 Hvis da $\sqrt{n}(m - M) \rightarrow N(0, \frac{1}{4} f'(M)^2)$
 der X har tetthet $f(x)$
 Men $f(x)$ er vanskelig å estimere
 $\rightarrow f'(x)$ enda vanskeligere.

sep 12-13:26

Hvorfor ikke alltid bootstrap?
 • Numerisk tungt
 • Lite ubehagelig at vi ikke får samme resultat-ekskald - ved gjentatt bootstrapping
 • For mange problemstillinger fungerer "klassiske" metoder godt nok.
 • Ikke alltid opplagt hvordan man skal bootstrappe, f.eks. med avhengige data.
 R kan imidlertid en automatisk farge for bootstrapping: boot
 Man må skrive litt mer (boot)
 først.

sep 12-13:49