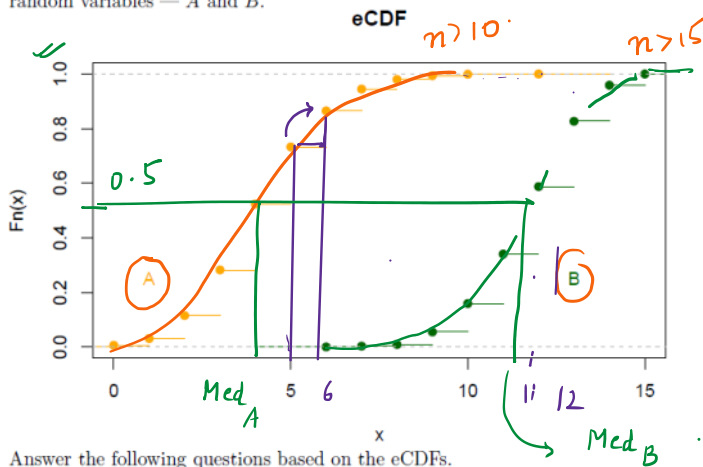


Below you can see a plot of two empirical cumulative distribution functions (eCDF). Each of the eCDFs are plotted using 1 000 realizations of one of two random variables — A and B .



$$A \sim \text{Bin}(n, p_A)$$

$$B \sim \text{Bin}(n, p_B)$$

Eg: A

cdf:

$$P[A \leq a] = \sum_{i=0}^a f(x)$$

Answer the following questions based on the eCDFs.

- i. Suppose both A and B are Binomial distributions with a common n and different p s. What is the value of the common parameter n ? Explain in one or two sentences.

Solution

Since both have the same parameter n , I can take the larger value for which there's a jump in the CDF. So $n = 15$.

- ii. Which of the 2 variables has a larger median? Explain in one or two sentences.

Solution

B has a larger median since its eCDF crosses 0.5 on the y-axis along at a greater point along the x-axis.

Question 2

A skate rental shop records the time between skates being rented and being returned. Analyzing 100 returns, the average time to return is found to be 2 hours. The shop opens for 4 hours daily and overnight rentals are not allowed.

Part a)

Suggest a reasonable parametric model among the models listed in Table 1 for the rental times assuming they are a random sample. What additional assumptions are you making with the selected distribution? Are any of the model assumptions unrealistic? Explain in a few sentences.

Solution

Given the time until rental return data, an Exponential distribution would be the most suitable choice among the distributions listed.

In addition to the each return time being independent and identical, the model assumes that the rental return is a Poisson process. That is, the likelihood of rental return at any moment is independent and identical.

The rental shop is open for 4 hours per day only. But an exponential random variable can take any positive real number.

Alternatively, you may have assumed a uniform distribution since the rental time is bounded. However, this is not appropriate since not every one rents their skates when the shop opens.

In either case, the independence assumption may be unrealistic since skaters rent together with friends and family. They are likely to return at the same time.

Part b)

What is your best guess for the parameter(s) of your selected model based on the given information? State any probability rules that support your guess.

Solution

Let X be the skate return time. If we assume $X \sim \text{Exp}(\lambda)$, $E(X) = \frac{1}{\lambda}$. The sample mean over the $n = 100$ repairs is $\bar{x}_{50} = 2$ hours. Thus, a reasonable estimate of $\lambda = 1/2 = 0.5$ based on the law of large numbers. (MOM/MLE)

X : Time for which skates are taken on rent.

$n = 100$

M.S X_1, X_2, \dots, X_{100}

$E(X) = 2, X \leq 4$

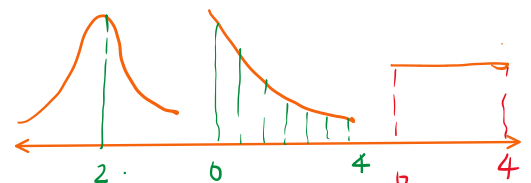
or $P[X > 4] = 0$

Range of $0 \leq X \leq 4$



Continuous variable

$N(\mu, \sigma^2)$ $\text{Exp}(\lambda)$ $U[a, b]$



\therefore M.S: $X_1, X_2, \dots, X_n \stackrel{i.i.d}{\sim} f(x) \Rightarrow X_i \sim f(x) \forall i$ and cdf $F(x)$

ordered sample: $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$

$\min\{X_1, X_2, \dots, X_n\}$

$\max\{X_1, X_2, \dots, X_n\}$

Q. Pdf of $X_{(n)} = \max\{X_1, X_2, \dots, X_n\}$

c.d.f of $X_{(n)} = P[X_{(n)} \leq x]$

$= P[\max\{X_1, X_2, \dots, X_n\} \leq x]$

$= P[X_1 \leq x, X_2 \leq x, \dots, X_n \leq x]$

$= P[X_1 \leq x] \cdot P[X_2 \leq x] \cdot \dots \cdot P[X_n \leq x]$

$\frac{d}{dx} F(x) = f(x)$

$$= P[X_1 \leq x] \cdot P[X_2 \leq x] \cdots P[X_n \leq x]$$

\downarrow \downarrow
 cdf of r.v X_1 cdf of r.v X_2
 $\simeq F(x)$ $\simeq F(x)$

$$= \underbrace{F(x) \cdot F(x) \cdots F(x)}_{n \text{ times}} = \{F(x)\}^n$$

pdf of $X_{(n)} = \frac{d}{dx} \{F(x)\}^n = n \{F(x)\}^{n-1} \cdot f(x)$.

g. pdf of $X_{(1)} = \min\{X_1, X_2, \dots, X_n\}$

cdf of $X_{(1)} = P[X_{(1)} \leq x]$

$$= P[\min\{X_1, X_2, \dots, X_n\} \leq x]$$

$$= 1 - P[\min\{X_1, X_2, \dots, X_n\} > x]$$

$$= 1 - P[X_1 > x, X_2 > x, \dots, X_n > x]$$

$$= 1 - \underbrace{P[X_1 > x]}_{\downarrow} \underbrace{P[X_2 > x]}_{\downarrow} \cdots P[X_n > x]$$

\downarrow \downarrow
 $1 - P[X_1 \leq x]$ $1 - P[X_2 \leq x]$
 $= \{1 - F(x)\}$ $= \{1 - F(x)\}$

$$= 1 - \{1 - F(x)\}^n$$

pdf of $X_{(1)} = \frac{d}{dx} [1 - \{1 - F(x)\}^n]$

$$= 0 - n \{1 - F(x)\}^{n-1} \cdot (0 - f(x))$$

$$= n \cdot f(x) \{1 - F(x)\}^{n-1}$$

$y = f(x)$
 $= \sin^n x$
 $= (\sin x)^n$
 $\frac{dy}{dx} = n \cdot \sin^{n-1} x \cdot \cos x$

$\sim \dots \dots \dots \{ \dots \dots \dots \}$

$$\frac{dy}{dz} = n \cdot \sin z^{n-1}$$

Question 3

Suppose $X_1, X_2, \dots, X_n \stackrel{i.i.d}{\sim} \text{Exp}(\lambda)$. Is the estimator $\hat{\lambda} = X_{(1)} = \min(X_1, X_2, \dots, X_n)$ an unbiased estimator?

Hint: \rightarrow cdf of $X_{(1)}$

$$P(X_{(1)} \leq x) = 1 - P(X_{(1)} > x) = 1 - P((X_1 > x) \cap (X_2 > x) \cap \dots \cap (X_n > x))$$

Solution:

$$\begin{aligned} P(X_{(1)} \leq x) &= 1 - P(X_{(1)} > x) = 1 - P((X_1 > x) \cap (X_2 > x) \cap \dots \cap (X_n > x)) \\ &= 1 - \prod_{i=1}^n [1 - P(X_i \leq x)] \\ &= 1 - \prod_{i=1}^n e^{-\lambda x} \\ &= 1 - e^{-n\lambda x} \text{ for } x \geq 0 \end{aligned}$$

$$\Rightarrow \text{cdf of } X_{(1)} = 1 - e^{-n\lambda x}$$

$$\Rightarrow X_{(1)} \sim \text{Exp}(n\lambda)$$

$$\Rightarrow E(\hat{\lambda}) = E(X_{(1)}) = \frac{1}{n\lambda}$$

Therefore $\hat{\lambda}$ is a biased estimator for λ .

$$\begin{aligned} \text{pdf of } X_{(1)} &= \frac{d}{dx} [1 - e^{-n\lambda x}] \\ &= -e^{-n\lambda x} (-n\lambda) \\ &= n\lambda e^{-n\lambda x} \end{aligned}$$

$$E[X_{(1)}] = \int_0^{\infty} x \cdot n\lambda e^{-n\lambda x} dx$$

Let $n\lambda x = z \Rightarrow x=0, z=0$

$$dx = \frac{dz}{n\lambda} \quad x \rightarrow \infty, z \rightarrow \infty$$

$$= \int_0^{\infty} n\lambda \cdot \left(\frac{z}{n\lambda}\right) \cdot e^{-z} \cdot \frac{dz}{n\lambda}$$

$$= \frac{1}{n\lambda} \left(\int_0^{\infty} z e^{-z} dz \right) = \frac{1}{n\lambda}$$

$$\int z e^{-z} dz = z \int e^{-z} dz = \int -e^{-z} dz$$

$$= z \int e^{-z} dz + \int e^{-z} dz$$

$$= -z e^{-z} - e^{-z} + C = -e^{-z} (z+1) + C$$

$$\int_0^{\infty} z e^{-z} dz = - \left[e^{-z} (z+1) \right]_0^{\infty} = - \left[e^{-\infty} (\infty+1) - e^{-0} (0+1) \right]$$

$$= - \left[0 - 1 \right] = 1$$

Bootstrap :-

$X \sim N(\mu, \sigma^2) \Rightarrow$ popln.

H.S: $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2) \rightarrow$ Theory

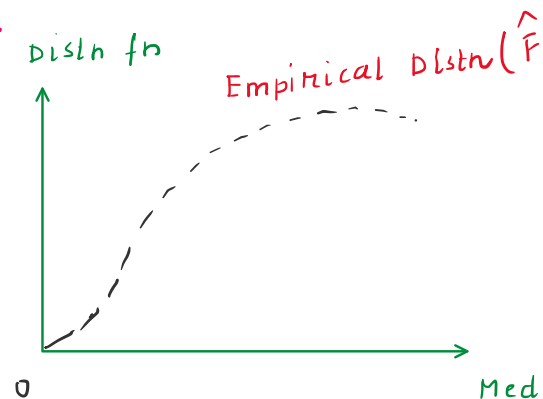
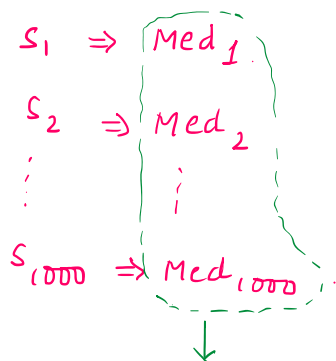
Sample mean: $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

Sample median: need not follow any particular known distribution

Bootstrap technique is used when the estimators do not follow a known stat distribution. Hence its properties cannot be analyzed directly from the popln distribution.

Obj: sampling distribution of the Med:

Take $n = 1000$ samples from popln.



Ordered set of med values: $Med_{(1)} \leq Med_{(2)} \leq \dots \leq Med_{(1000)}$

Q. $Var[\text{sample}(Med)] \Rightarrow$ obtained from \hat{F} .

$$\text{Mean}[\text{sample}(Med)] = \bar{x}_{Med} = \frac{1}{1000} \sum_{i=1}^{1000} Med_{(i)}$$

$$\text{Mean [sample (Med)]} = \bar{x}_{\text{Med}} = \frac{1}{1000} \sum_{i=1}^{1000} \text{Med}_{(i)}$$

$$\text{Var [sample (Med)]} = \frac{1}{1000-1} \sum_{i=1}^n (\text{Med}_{(i)} - \bar{x}_{\text{Med}})^2$$