

STAT/MTHE 353: 2 – Order Statistics

Order Statistics

Let X_1, \dots, X_n be jointly continuous random variables.

Definition The k th order statistic of X_1, \dots, X_n is the k th smallest value of the X_i and denoted by $X_{(k)}$. Thus

$$X_{(1)} = \min(X_1, \dots, X_n), \quad X_{(n)} = \max(X_1, \dots, X_n),$$

The vector $(X_{(1)}, \dots, X_{(n)})$ is called the order statistics of (X_1, \dots, X_n) .

- We will assume that X_1, \dots, X_n are *independent and identically distributed* (i.i.d.) random variables. Thus their joint pdf is

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = f(x_1)f(x_1) \cdots f(x_n)$$

where f is the *common* marginal pdf of the X_i .

- We will determine the joint pdf of $X_{(1)}, \dots, X_{(n)}$.

If $f_{1\dots n}$ is the joint pdf of $X_{(1)}, \dots, X_{(n)}$, then (with $h > 0$)

$$\begin{aligned} f_{1\dots n}(x_1, \dots, x_n) &= \lim_{h \rightarrow 0} \frac{1}{h^n} \int_{x_1}^{x_1+h} \cdots \int_{x_n}^{x_n+h} f_{1\dots n}(t_1, \dots, t_n) dt_1 \cdots dt_n \\ &= \lim_{h \rightarrow 0} \frac{1}{h^n} P(x_1 \leq X_{(1)} \leq x_1 + h, \dots, x_n \leq X_{(n)} \leq x_n + h) \end{aligned}$$

(true for any sufficiently well behaved pdf).

Clearly, we only need consider the case $x_1 < x_2 < \cdots < x_n$. (Why?) Let S_n denote the set of permutations of $\{1, \dots, n\}$. Note that $|S_n| = n!$.

Since X_1, \dots, X_n are jointly continuous, the $n!$ *disjoint* events

$$A_\sigma = \{X_{\sigma(1)} < X_{\sigma(2)} < \cdots < X_{\sigma(n)}\}, \quad \sigma \in S_n$$

have total probability 1, i.e.,

$$P\left(\bigcup_{\sigma \in S_n} A_\sigma\right) = 1.$$

Define the events

$$A = \{x_1 \leq X_{(1)} \leq x_1 + h, \dots, x_n \leq X_{(n)} \leq x_n + h\}$$

and

$$B_\sigma = \{x_1 \leq X_{\sigma(1)} \leq x_1 + h, \dots, x_n \leq X_{\sigma(n)} \leq x_n + h\}$$

and note that $A \cap A_\sigma = B_\sigma$ if $h > 0$ is small enough. Thus for such h

$$\begin{aligned} P(x_1 \leq X_{(1)} \leq x_1 + h, \dots, x_n \leq X_{(n)} \leq x_n + h) &= P(A) = \sum_{\sigma \in S_n} P(A \cap A_\sigma) = \sum_{\sigma \in S_n} P(B_\sigma) \\ &= \sum_{\sigma \in S_n} \prod_{i=1}^n P(x_i \leq X_{\sigma(i)} \leq x_i + h) \quad (\text{since the } X_i \text{ are independent}) \\ &= \sum_{\sigma \in S_n} \prod_{i=1}^n [F_{X_{\sigma(i)}}(x_i + h) - F_{X_{\sigma(i)}}(x_i)] \\ &= n! \prod_{i=1}^n [F(x_i + h) - F(x_i)] \quad (F \text{ is the common cdf of the } X_i) \end{aligned}$$

We obtain

$$\begin{aligned}
 f_{1\dots n}(x_1, \dots, x_n) &= \lim_{h \rightarrow 0} \frac{1}{h^n} P(x_1 \leq X_{(1)} \leq x_1 + h, \dots, x_n \leq X_{(n)} \leq x_n + h) \\
 &= \lim_{h \rightarrow 0} \frac{1}{h^n} n! \prod_{i=1}^n [F(x_i + h) - F(x_i)] \\
 &= n! \prod_{i=1}^n \lim_{h \rightarrow 0} \frac{F(x_i + h) - F(x_i)}{h} \\
 &= n! f(x_1) f(x_2) \cdots f(x_n).
 \end{aligned}$$

In conclusion, the joint pdf of the order statistics $X_{(1)}, \dots, X_{(n)}$ is given by

$$f_{1\dots n}(x_1, \dots, x_n) = \begin{cases} n! f(x_1) f(x_2) \cdots f(x_n) & \text{if } x_1 < x_2 < \cdots < x_n, \\ 0 & \text{otherwise.} \end{cases}$$

Marginal pdfs of order statistics

Notation: For $k < r$ we let $f_{k\dots r}(x_k \cdots x_r)$ denote the marginal jpdf of $X_{(k)}, \dots, X_{(r)}$. Here we always assume $x_k < x_{k+1} < \cdots < x_r$; otherwise $f_{k\dots r}(x_k \cdots x_r) = 0$.

“Integrating out” x_1 we get

$$\begin{aligned}
 f_{2\dots n}(x_2, \dots, x_n) &= \int_{-\infty}^{x_2} n! f(x_1) f(x_2) \cdots f(x_n) dx_1 \\
 &= n! F(x_2) f(x_2) \cdots f(x_n)
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 f_{3\dots n}(x_3, \dots, x_n) &= \int_{-\infty}^{x_3} n! F(x_2) f(x_2) \cdots f(x_n) dx_2 \\
 &= n! \frac{F(x_3)^2}{2} f(x_3) \cdots f(x_n)
 \end{aligned}$$

Integrating out x_3, \dots, x_{k-1} , we obtain in general for all $x_k < \cdots < x_n$

$$f_{k\dots n}(x_k, \dots, x_n) = n! \frac{F(x_k)^{k-1}}{(k-1)!} f(x_k) \cdots f(x_n)$$

If we integrate out x_n first, we obtain

$$\begin{aligned}
 f_{1\dots n-1}(x_1, \dots, x_{n-1}) &= \int_{x_{n-1}}^{\infty} n! f(x_1) f(x_2) \cdots f(x_n) dx_n \\
 &= n! f(x_1) \cdots f(x_{n-1}) (1 - F(x_{n-1}))
 \end{aligned}$$

Continuing with $x_{n-1}, x_{n-2}, \dots, x_{r+1}$, we get for all $x_1 < \cdots < x_r$

$$f_{1\dots r}(x_1, \dots, x_r) = n! f(x_1) \cdots f(x_r) \frac{(1 - F(x_r))^{n-r}}{(n-r)!}$$

Integrating out $x_n, x_{n-1}, \dots, x_{r+1}$ in $f_{k\dots n}(x_k, \dots, x_n)$ or x_1, x_2, \dots, x_{k-1} in $f_{1\dots r}(x_1, \dots, x_r)$ we finally obtain for all $x_k < \cdots < x_r$

$$f_{k\dots r}(x_k, \dots, x_r) = n! \frac{F(x_k)^{k-1}}{(k-1)!} f(x_k) \cdots f(x_r) \frac{(1 - F(x_r))^{n-r}}{(n-r)!}$$

Setting $r = k + 1$ we get for all $x_k < x_{k+1}$

$$f_{k,k+1}(x_k, x_{k+1}) = n! \frac{F(x_k)^{k-1}}{(k-1)!} f(x_k) f(x_{k+1}) \frac{(1 - F(x_{k+1}))^{n-k-1}}{(n-k-1)!}$$

From this we obtain the marginal pdf of $X_{(k)}$:

$$\begin{aligned}
 f_k(x_k) &= n! \frac{F(x_k)^{k-1}}{(k-1)!} f(x_k) \int_{x_k}^{\infty} f(x_{k+1}) \frac{(1 - F(x_{k+1}))^{n-k-1}}{(n-k-1)!} dx_{k+1} \\
 &= n! \frac{F(x_k)^{k-1}}{(k-1)!} f(x_k) \left[-\frac{(1 - F(x_{k+1}))^{n-k}}{(n-k)!} \right]_{x_k}^{\infty} \\
 &= n! \frac{F(x_k)^{k-1}}{(k-1)!} f(x_k) \frac{(1 - F(x_k))^{n-k}}{(n-k)!}
 \end{aligned}$$

We obtained (Theorem 9.5 in text):

$$f_k(x) = \frac{n!}{(k-1)!(n-k)!} f(x) F(x)^{k-1} (1 - F(x))^{n-k}$$

Finally, if we integrate out x_{k+1}, \dots, x_{r-1} in $f_{k \dots r}$, then we obtain for all $x_k < x_r$ (Theorem 9.6):

$$f_{k,r}(x_k, x_r) = n! \frac{f(x_k)F(x_k)^{k-1}}{(k-1)!} \cdot \frac{(F(x_r) - F(x_k))^{r-k-1}}{(r-k-1)!} \cdot \frac{f(x_r)(1 - F(x_r))^{n-r}}{(n-r)!}$$

Special cases: minimum and maximum

From f_k with $k = 1$ we get the pdf of $X_{(1)} = \min(X_1, \dots, X_n)$:

$$f_1(x) = n f(x) (1 - F(x))^{n-1}$$

For $k = n$ we get the pdf of $X_{(n)} = \max(X_1, \dots, X_n)$:

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

Example: The i.i.d. random variables X_1, X_2, \dots, X_n are often called a *random sample* from the common distribution of the X_i . The *range* of the random sample is $R = X_{(n)} - X_{(1)}$. Determine the pdf of R . Specialize to the case where $X_i \sim \text{Uniform}(0, 1)$.

Solution: ...

Application: Sample median

- The *median* of the distribution of a r.v. X is any value m such that

$$P(X \leq m) \geq \frac{1}{2} \quad \text{and} \quad P(X \geq m) \geq \frac{1}{2}.$$

Equivalently, if F is the cdf of X ,

$$F(m) \geq \frac{1}{2} \quad \text{and} \quad 1 - F(m^-) \geq \frac{1}{2}.$$

- The *sample median* of the random sample X_1, \dots, X_n from a continuous distribution is

$$M = \begin{cases} X_{(\frac{n+1}{2})} & \text{if } n \text{ is odd} \\ \frac{X_{(\frac{n}{2})} + X_{(\frac{n+1}{2})}}{2} & \text{if } n \text{ is even} \end{cases}$$

- The sample median is often taken to be an estimate of the median, and it is sometimes preferred to the *sample mean* $\frac{1}{n} \sum_{i=1}^n X_i$ as an estimate of the “center of the distribution” because it is more robust to “outliers.”

Example: ...