

Statistical Bioinformatics, Makerere Extreme Values & Order Statistics Timo Koski

TK

02.08.2018



The lecture will cover Chapter 1 and 2.11 in Ewens and Grant in some detail.



A Extreme values

- 1) Minima and maxima of n I.I.D. RV's
- 2) Memoryless property, mean and variance of maximum of n exponential R.V.'s
- 3) Gumbel's distribution: an asymptotic extreme value distribution

B Order statistics

- 4) Order Statistics: definition, joint distribution
- 5) Probabilities for pieces of a randomly broken stick

Minimum of n I.I.D. R.V's

X_1, X_2, \dots, X_n , I.I.D with the same distribution as X . Let

$$X_{\min} \stackrel{\text{def}}{=} \min (X_1, X_2, \dots, X_n).$$

Then

$$P(X_{\min} \geq x) = (P(X \geq x))^n.$$



Proof:

$$\begin{aligned} P(X_{\min} \geq x) &= P(\min(X_1, X_2, \dots, X_n) \geq x) = \\ &= P(X_1 \geq x, X_2 \geq x, \dots, X_n \geq x) = \prod_{i=1}^n P(X_i \geq x), \end{aligned}$$

by independence of X_1, X_2, \dots, X_n .



Minimum of n I.I. Exponential R.V's

Recall that X_1, X_2, \dots, X_n , I.I. and $\in \text{Exp}(\lambda)$

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & 0 < x < \infty \\ 0 & \text{elsewhere.} \end{cases}$$

$$F_X(x) = 1 - e^{-\lambda x}$$

$$P(X \geq x) = 1 - F_X(x) = 1 - (1 - e^{-\lambda x}) = e^{-\lambda x}.$$



Minimum of n I.I. Exponential R.V.'s

X_1, X_2, \dots, X_n , I.I. and $\text{Exp}(\lambda)$. $X_{\min} \stackrel{\text{def}}{=} \min(X_1, X_2, \dots, X_n)$.

$$\begin{aligned} P(X_{\min} \geq x) &= (P(X \geq x))^n = \\ &= (e^{-\lambda x})^n = e^{-n\lambda x} \end{aligned}$$

Hence $X_{\min} \in \text{Exp}(n\lambda)$.



Maximum of n I.I.D. R.V's

X_1, X_2, \dots, X_n , I.I.D with the same distribution as X . Let

$$X_{\max} \stackrel{\text{def}}{=} \max(X_1, X_2, \dots, X_n).$$

$$P(X_{\max} \leq x) = (P(X \leq x))^n$$



Maximum of n I.I.D. R.V's

$$\begin{aligned} P(X_{\max} \leq x) &= P(\max(X_1, X_2, \dots, X_n) \leq x) = \\ &= P(X_1 \leq x, X_2 \leq x, \dots, X_n \leq x) = \prod_{i=1}^n P(X_i \leq x), \end{aligned}$$

by independence of X_1, X_2, \dots, X_n . □



Maximum of n I.I.D. Exponential R.V's

X_1, X_2, \dots, X_n , i.i. and $\text{Exp}(\lambda)$. $X_{\max} \stackrel{\text{def}}{=} \max(X_1, X_2, \dots, X_n)$.

$$\begin{aligned} P(X_{\max} \leq x) &= (P(X \leq x))^n = \\ &= (1 - e^{-\lambda x})^n. \end{aligned}$$



Now we shall compute

$$E [X_{\max}], \text{Var} [X_{\max}]$$

and an asymptotic distribution of X_{\max} . This will be a lengthier story, some details will be skipped. We need the memoryless property of the exponential distributions. We prove this for the geometric distribution.

Memoryless property: Geometric distribution

Recall $X \in Ge(p)$

$$f_X(x) = (1-p) \cdot p^x; x = 0, 1, \dots,$$

We seek the conditional probability

$$P(X \geq x_1 + x_2 \mid X \geq x_1).$$

By definition

$$P(X \geq x_1 + x_2 \mid X \geq x_1) = \frac{P(X \geq x_1 + x_2)}{P(X \geq x_1)}.$$

We need the distribution function of $X \in Ge(p)$.



Distribution function of $X \in \text{Ge}(p)$

$$\begin{aligned}F_X(x) &= P(X \leq x) = \sum_{k=0}^x (1-p) \cdot p^k = \\&= (1-p) \cdot \sum_{k=0}^x p^k = (1-p) \frac{1-p^{x+1}}{(1-p)} \\&= 1-p^{x+1},\end{aligned}$$

by summing a geometric sum.



Memoryless property: Geometric distribution

Then

$$P(X \geq x_1 + x_2 \mid X \geq x_1) = \frac{P(X \geq x_1 + x_2)}{P(X \geq x_1)} = \frac{p^{x_1+x_2+1}}{p^{x_1+1}} = p^{x_2},$$

since $P(X \leq x) = 1 - p^{x+1}$. This is the memoryless property.



$X \in \text{Exp}(\lambda)$. We find as above the conditional probability

$$P(X \geq x_1 + x_2 \mid X \geq x_1) = \frac{P(X \geq x_1 + x_2)}{P(X \geq x_1)} = e^{-\lambda x_2}$$

$$\frac{d}{ds}\phi(0) = E[X]$$

$X \in \text{Exp}(\lambda)$

$$\phi_X(s) = \frac{\lambda}{\lambda - s}$$

\Rightarrow

$$\frac{d}{ds}\phi_X(s) = -\frac{-\lambda}{[(\lambda - s)]^2}$$

Hence $E[X] = \frac{1}{\lambda}$.

Mean and Variance of Max. of n Exponential R.V's

Now we get to computing

$$E[X_{\max}], \text{Var}[X_{\max}]$$

We are going to use the fact that $X_{\min} \in \text{Exp}(n\lambda)$ and the memoryless property.



Mean and Variance of Max. of n Exponential R.V's

Suppose, e.g., that n cellular proteins have lifetimes until degradation that are exponentially distributed. Then the expected time to one or other protein dies is the mean of the minimum,

$$E(X_{\min}) = \frac{1}{n\lambda},$$

since $X_{\min} \in \text{Exp}(n\lambda)$. By the memoryless property we have that the mean time till next protein dies is the minimum of $(n-1)$ exponential r.v.'s independent of the first death,

$$\frac{1}{(n-1)\lambda}$$

and so on.



Mean and Variance of Max. of n Exponential R.V's

The expected time till the last of the n proteins dies is the expectation of the maximum lifetime and is thus

$$\frac{1}{\lambda} + \frac{1}{2\lambda} + \cdots + \frac{1}{n\lambda} = E[X_{\max}].$$

Similarly we can find

$$\text{Var}[X_{\max}] = \frac{1}{\lambda^2} + \frac{1}{(2\lambda)^2} + \cdots + \frac{1}{(n\lambda)^2}.$$

Euler's constant

We write

$$\frac{1}{\lambda} + \frac{1}{2\lambda} + \cdots + \frac{1}{n\lambda} = \frac{1}{\lambda} \sum_{i=1}^n \frac{1}{i}$$

It is shown in courses in mathematics that

$$\lim_{n \rightarrow \infty} \left(\sum_{i=1}^n \frac{1}{i} - \log n \right) = \gamma,$$

where γ is Euler's constant, (see Westergren & Råde: BETA)



Mean and Variance of Max. of n Exponential R.V's

Hence we approximate

$$E[X_{\max}] \approx \frac{\gamma + \log n}{\lambda}.$$

\log is the natural logarithm.

$$\text{Var}[X_{\max}] \approx \frac{\pi^2}{6\lambda^2}.$$



Asymptotic Extreme Value Distribution (EVD) (1.1)

$X_{\max} = \max(X_1, X_2, \dots, X_n)$, X_1, X_2, \dots, X_n , I.I.D.

Note that X_1, X_2, \dots, X_n need not be exponentially distributed, but the distribution must satisfy a tail condition (see below). We are going to prove

$$\lim_{n \rightarrow \infty} P\left(X_{\max} - \frac{\log n}{\lambda} \leq x\right) = e^{-e^{-\lambda x}}.$$

Here log is the natural logarithm.



Asymptotic EVD (1.2): Gumbel's distribution

$$F(x) := e^{-e^{-\lambda x}}$$

is a probability distribution function:

- $F(x) \rightarrow 1$, as $x \rightarrow \infty$,
- $F(x) \rightarrow 0$, as $x \rightarrow -\infty$,
- $F(x)$ is an increasing function, since $\frac{d}{dx}F(x) = \lambda e^{-\lambda x - e^{-\lambda x}} > 0$

The distribution is known as Gumbel's distribution.



$$F_n(x) = P\left(X_{\max} - \frac{\log n}{\lambda} \leq x\right).$$

We claim thus

$$\lim_{n \rightarrow \infty} F_n(x) = F(x).$$

We say that X_{\max} converges asymptotically in distribution to a random variable that has a Gumbel distribution.

In general, if $F_1(x), F_2(x), F_3(x), F_4(x), \dots$ are distribution functions, and

$$\lim_{n \rightarrow \infty} F_n(x) = F(x).$$

where $F(x)$ is a distribution function, we talk about convergence in distribution. Examples: Central Limit Theorem.

EVD (2): Proof starts

We know from the preceding that

$$\begin{aligned} P\left(X_{\max} - \frac{\log n}{\lambda} \leq x\right) &= P\left(X_{\max} \leq x + \frac{\log n}{\lambda}\right) \\ &= \left[P\left(X_1 \leq x + \frac{\log n}{\lambda}\right)\right]^n \\ &= \left[F\left(x + \frac{\log n}{\lambda}\right)\right]^n. \end{aligned}$$

where $F(x) = P(X_1 \leq x)$.

We write the right hand side as

$$\left[F \left(x + \frac{\log n}{\lambda} \right) \right]^n = e^{n \log F \left(x + \frac{\log n}{\lambda} \right)},$$

EVD (4): Assumption on tails

We assume that

$$F(x) = 1 - V(x)e^{-\lambda x}$$

for large x , where $\lim_{x \rightarrow \infty} V(x) = 1$ is assumed for simplicity. In words, the tail of the distribution function

$$1 - F(x) = V(x)e^{-\lambda x}$$

vanishes exponentially. Clearly the exponential distribution satisfies this !



EVD (5): Insert the tail

$$e^{n \log F\left(x + \frac{\log n}{\lambda}\right)} = e^{n \log \left[1 - V\left(x + \frac{\log n}{\lambda}\right) e^{-x\lambda - \log n}\right]}.$$



EVD (6): Series expansion

We use the series expansion (valid for $-1 \leq x < 1$)

$$\log(1+x) = x - \frac{1}{2}x^2 + \dots$$

with the result

$$\begin{aligned} n \log \left[1 - V \left(x + \frac{\log n}{\lambda} \right) e^{-x\lambda} e^{\log 1/n} \right] &= \\ &= -n \frac{1}{n} V \left(x + \frac{\log n}{\lambda} \right) e^{-x\lambda} - \frac{1}{2} n \frac{1}{n^2} G(n) + O \left(\frac{1}{n} \right) \dots \end{aligned}$$

where $G(n)$ is a bounded function of n and $O\left(\frac{1}{n}\right)$ represents the omitted terms that turn to zero as $n \rightarrow +\infty$.

Hence we have

$$\begin{aligned}\lim_{n \rightarrow \infty} n \log \left[1 - V \left(x + \frac{\log n}{\lambda} \right) e^{-x\lambda} e^{\log 1/n} \right] \\ = -e^{-x\lambda}.\end{aligned}$$

In other words we have found that

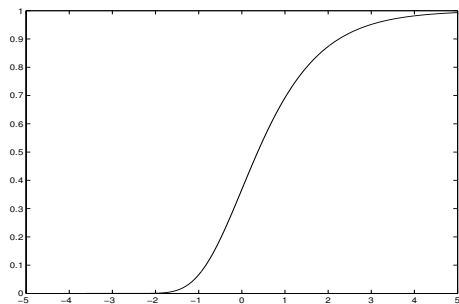
$$\lim_{n \rightarrow \infty} P \left(X_{\max} - \frac{\log n}{\lambda} \leq x \right) = e^{-e^{-\lambda x}},$$

as was to be proved.

The rate of convergence may be slow.

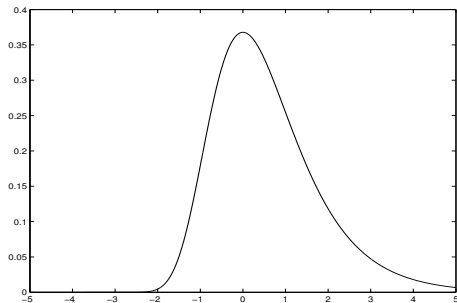
Gumbel's distribution

$$P(X_{\max} \leq x) \approx e^{-e^{-u}}$$



Gumbel density

$$e^{-u} - e^{-e^{-u}}$$



The p -value of an observed value x of X_{\max} is the probability $P(X_{\max} > x)$ under some model for the underlying random variables. If these are exponential, then

$$p\text{-value} \approx 1 - e^{-ne^{-\lambda x}}.$$

The smaller the p -value, the less likely is the hypothesis that the model is the correct one.

B: Order Statistics: the setting

Let X_1, \dots, X_n be I.I.D. random variables distributed like X with the distribution function F such that

$$F(x) = \int_{-\infty}^x f(u) du.$$

We are for reasons of simplicity of writing going to drop here the subscripts in $F_X(x) = \int_{-\infty}^x f_X(u) du$.



X_1, \dots, X_n are I.I.D. random variables distributed like X .

$X_{(j)}$ = j :th order statistic =

j :th smallest of X_1, \dots, X_n .

$X_{(1)} = X_{\min}$, $X_{(n)} = X_{\max}$. Clearly

$$X_{(1)} < X_{(2)} < \dots < X_{(n-1)} < X_{(n)}$$

Order Statistics: definition (remark)

X_1, \dots, X_n are I.I.D. random variables distributed like X . Since X is a continuous r.v. , we have strict equalities in

$$X_{(1)} < X_{(2)} < \dots < X_{(n-1)} < X_{(n)}$$

For a continuous random variable $P(X = x) = \int_x^x f(u) du = 0$.



Order Statistics: the theorem on the joint density

X_1, \dots, X_n are I.I.D. random variables distributed like X with density f .
Then

$$f_{X_{(1)}, \dots, X_{(n)}}(x_{(1)}, \dots, x_{(n)}) = \begin{cases} n! \prod_{i=1}^n f(x_{(i)}) & \text{if } x_{(1)} < \dots < x_{(n)} \\ 0 & \text{elsewhere.} \end{cases} \quad \square$$

$$f_{X_{(1)}, \dots, X_{(n)}}(x_{(1)}, \dots, x_{(n)}) \\ = \begin{cases} n! \prod_{i=1}^n f(x_{(i)}) & \text{if } x_{(1)} < \dots < x_{(n)} \\ 0 & \text{elsewhere.} \end{cases} \quad \square$$

The theorem will be explained on the whiteboard.

A Tool: Multinomial Probability Distribution

The probability

$$P(Y_1 = n_1, Y_2 = n_2, \dots, Y_L = n_L) = \frac{n!}{n_1! n_2! \dots n_L!} \cdot p_1^{n_1} \cdots p_L^{n_L}.$$

is the joint probability of the random variables Y_1, Y_2, \dots, Y_L as functions of the random variables X_1, \dots, X_n by

$$Y_j = \text{the number of times } X_i = x_j, \quad i = 1, 2, \dots, n,$$

with L different outcomes x_j .



Let $0 < L$ and

$$f_X(x) = \begin{cases} \frac{1}{L} & 0 < x < L \\ 0 & \text{elsewhere.} \end{cases}$$

$X \in \mathcal{U}(0, L)$, uniformly distributed.

X_1, \dots, X_n are I.I.D. random variables $X \in U(0, L)$

$$f_{X_{(i)}}(x_{(i)}) = \frac{n!}{(i-1)!(n-i)!} x_{(i)}^{i-1} (L - x_{(i)})^{n-i} L^{-n}$$

X_1, \dots, X_n are I.I.D. random variables $X \in U(0, L)$

$$f_{X_{\min}}(x_{\min}) = f_{X_{(1)}}(x_{(1)}) = n \cdot (L - x_{\min})^{n-1} L^{-n}.$$

$X \in U(0, 1)$

$$f_{X_{(i)}}(x_{(i)}) = \frac{n!}{(i-1)!(n-i)!} x_{(i)}^{i-1} (1-x_{(i)})^{n-i}.$$

This is Beta($i, n - i + 1$).

X_1, \dots, X_n are I.I.D. random variables $X \in U(0, L)$

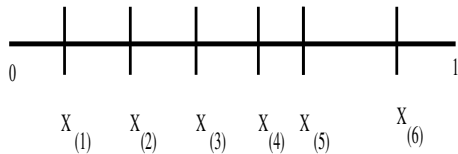
$$f_{X_{(1)}, \dots, X_{(n)}}(x_{(1)}, \dots, x_{(n)}) = \begin{cases} \frac{n!}{L^n} & \text{if } x_{(1)} < \dots < x_{(n)} \\ 0 & \text{elsewhere.} \end{cases} \quad \square$$

X_1, \dots, X_n are I.I.D. random variables $X \in U(0, L)$

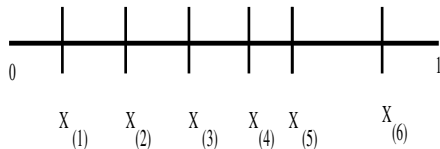
$$f_{X_1, \dots, X_n | X_{(1)}, \dots, X_{(n)}} \left(x_1, \dots, x_n \mid x_{(1)}, \dots, x_{(n)} \right) \\ = \begin{cases} \frac{1}{n!} & \text{if } x_{(1)} < \dots < x_{(n)} \\ 0 & \text{elsewhere.} \end{cases} \quad \square$$

A Stick Broken at Random

X_1, \dots, X_n are i.i.d. random variables $X \in U(0, 1)$.



A Stick Broken at Random

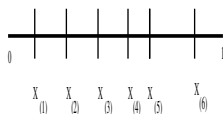


$$U_1 = X_{(1)} \quad U_2 = X_{(2)} - X_{(1)} \quad U_3 = X_{(3)} - X_{(2)} \quad U_4 = X_{(4)} - X_{(3)}$$

$$U_5 = X_{(5)} - X_{(4)} \quad U_6 = X_{(6)} - X_{(5)} \quad U_7 = 1 - X_{(6)}$$

U_1, \dots, U_{n+1} are the lengths of the $n + 1$ intervals created by X_1, \dots, X_n i.i.d. $\in U(0, 1)$.

A Stick Broken at Random



$$U_1 = X_{(1)} \quad U_2 = X_{(2)} - X_{(1)} \quad U_3 = X_{(3)} - X_{(2)} \quad U_4 = X_{(4)} - X_{(3)}$$

$$U_5 = X_{(5)} - X_{(4)} \quad U_6 = X_{(6)} - X_{(5)} \quad U_7 = 1 - X_{(6)}$$

We can show that

$$P(U_1 > u, \dots, U_g > u) = (1 - gu)^n$$

for if $1 - gu > 0$. Proof handled during a lektion.

Another result

Let $U_{(n-j)}$ be the $n - j$ th order statistic

$$P\left(U_{(n-j)} \leq x\right) \\ = \sum_{l=0}^j \binom{n+1}{l} \sum_{v=0}^{n+1-l} (-1)^v \binom{n+1-l}{v} (1 - (l+v)x)_+$$

where $(a)_+ = \max(0, a)$.

The exact distribution of U_{\max}

Set $j = 0$ in the above

$$P(U_{\max} \leq x) \\ = \sum_{v=0}^{n+1} (-1)^v \binom{n+1}{v} (1 - vx)_+$$

where $(a)_+ = \max(0, a)$.

Covering the whole genome

The probability of covering a whole genome by n 'markers' of length at most x is also equal to

$$P(U_{\max} \leq x) = \sum_{\nu=0}^{n+1} (-1)^\nu \binom{n+1}{\nu} (1 - \nu x)_+$$

where $(a)_+ = \max(0, a)$. K. Lange & M. Boehnke (1982): How many polymorphic genes will it take to span the human genome. *American Journal of Human Genetics*. 34, pp. 842–845.



Asymptotic distribution of U_{\max}

$$P(nU_{\max} - \ln n \leq x) = P(nU_n - \ln n \leq x) \rightarrow e^{-e^{-x}},$$

as $n \rightarrow \infty$. Found in (L. Holst 1980).



The number N_α of 'markers' of needed until they are at at most length α :

$$P(N_\alpha \leq n) = P(U_{\max} \leq \alpha) \approx e^{-e^{-x}}$$

if $n\alpha - \ln n \rightarrow x$, when $\alpha \rightarrow 0$, $n \rightarrow \infty$.

$e^{-x} = 10^{-c}$ gives $n \geq c \ln 10 / \alpha - \ln n / \alpha$.

Transformations of Variable (1)

$$U_i = U_i(X_1, \dots, X_n), \quad i = 1, \dots, n.$$

Then

$$f_{U_1, \dots, U_n}(u_1, \dots, u_n) = \\ f_{X_1, \dots, X_n}(x_1, \dots, x_n) \cdot |J^{-1}|,$$

J is known as the Jacobian of the transformation.

Transformations of Variable (1)

$$f_{U_1, \dots, U_n}(u_1, \dots, u_n) = f_{X_1, \dots, X_n}(x_1, \dots, x_n) \cdot |J^{-1}|,$$

where the Jacobian is

$$J = \begin{vmatrix} \frac{\partial U_1}{\partial X_1} & \frac{\partial U_1}{\partial X_2} & \cdots & \frac{\partial U_1}{\partial X_n} \\ \frac{\partial U_2}{\partial X_1} & \frac{\partial U_2}{\partial X_2} & \cdots & \frac{\partial U_2}{\partial X_n} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial U_n}{\partial X_1} & \frac{\partial U_n}{\partial X_2} & \cdots & \frac{\partial U_n}{\partial X_n} \end{vmatrix}.$$

This and an alternative formula are proved in courses on multivariable integral calculus.



Transformations of Variable (2)

$$X_i = X_i(U_1, \dots, U_n), \quad i = 1, \dots, n.$$

$$f_{U_1, \dots, U_n}(u_1, \dots, u_n) = \\ f_{X_1, \dots, X_n}(x_1, \dots, x_n) \cdot |J^*|,$$

where

$$J^* = \begin{vmatrix} \frac{\partial X_1}{\partial U_1} & \frac{\partial X_1}{\partial U_2} & \cdots & \frac{\partial X_1}{\partial U_n} \\ \frac{\partial X_2}{\partial U_1} & \frac{\partial X_2}{\partial U_2} & \cdots & \frac{\partial X_2}{\partial U_n} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial X_n}{\partial U_1} & \frac{\partial X_n}{\partial U_2} & \cdots & \frac{\partial X_n}{\partial U_n} \end{vmatrix}.$$

Proved in courses on multivariable integral calculus.

End of Lecture



KTH Matematik

