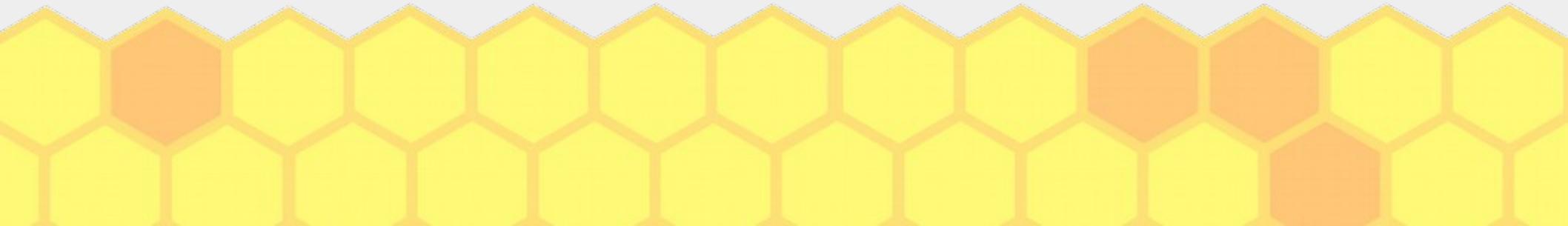


# Stochastic methods in Mathematical Modelling

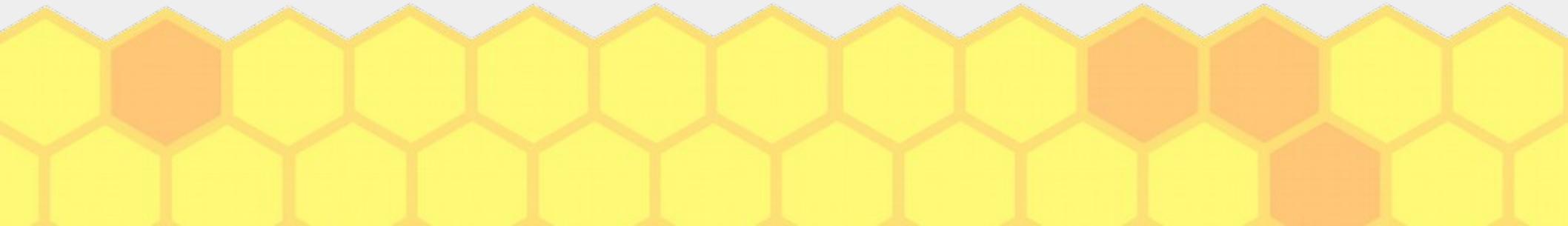
## Lecture 4. Extreme value statistics



# Extreme value statistics

## Applications

- Equity risks. Day to day market risk
- Extreme weather events
- Performances in a sport (what would be the next 100 m sprint record?)
- Mutational events during evolution
- Side effects of drugs  
etc. etc.



We are interested in the most extreme value  $M$  out of the set of random variables  $\{x_1, x_2, \dots, x_N\}$ .

$$\text{maximum } M = \max (x_1, x_2, \dots, x_N)$$

Simplification:  $x_1, \dots, x_n$  are i.i.d.

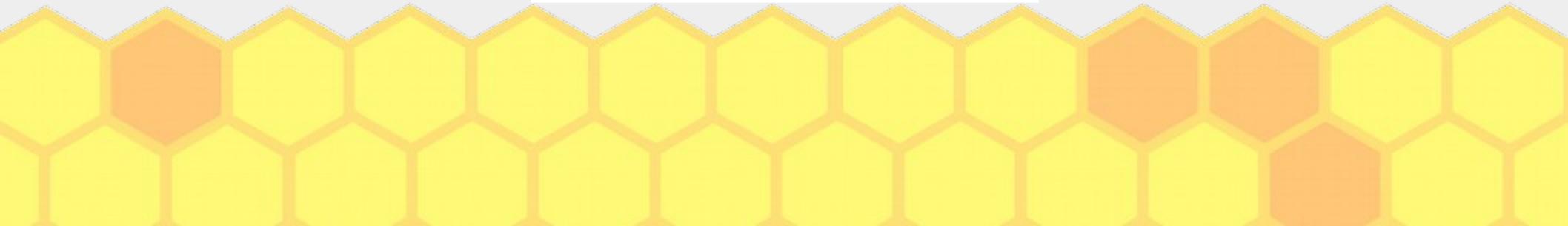
$$P(x_1, x_2, \dots, x_N) = p(x_1)p(x_2)\dots p(x_N)$$

For the mean of  $N$  samples we have a universality law (if the variance is finite) in form of CLT

What about the maximum?



*Universality in the case of i.i.d.*



We are interested in the most extreme value  $M$  out of the set of random variables  $\{x_1, x_2, \dots, x_N\}$ .

The general case CDF of distribution of  $M$  reads

$$\begin{aligned} Q_N(x) &= \text{Prob}[M \leq x, N] = \text{Prob}[x_1 \leq x, x_2 \leq x, \dots, x_N \leq x] , \\ &= \int_{-\infty}^x dx_1 \int_{-\infty}^x dx_2 \cdots \int_{-\infty}^x dx_N P(x_1, x_2, \dots, x_N) , \end{aligned}$$

Hence the pdf

$$P(M, N) = Q'_N(M) \equiv \frac{dQ_N}{dx}|_{x=M}$$



# Extreme value statistics

We are interested in the most extreme value  $M$  out of the set of random variables  $\{x_1, x_2, \dots, x_N\}$ .

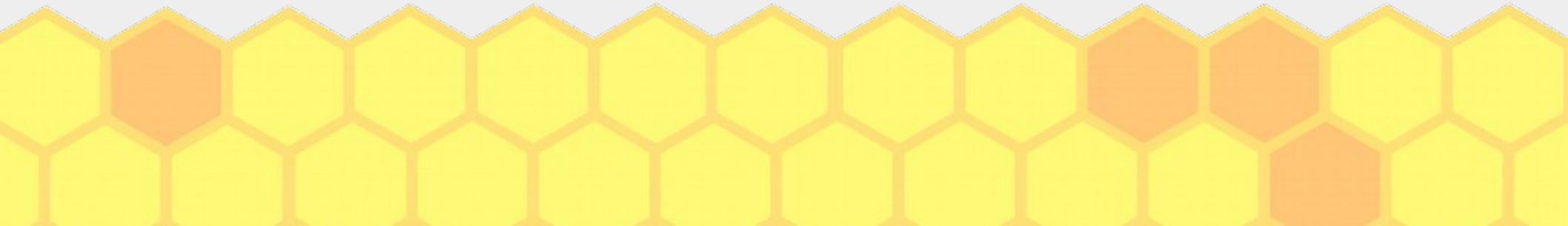
CDF for i.i.d. case

$$Q_N(x) = \left[ \int_{-\infty}^x dy \ p(y) \right]^N = \left[ 1 - \int_x^\infty dy \ p(y) \right]^N$$

When  $N$  tends to  $\infty$   $Q_N(x)$  approaches the limiting form

$$Q_N(x) \xrightarrow[\substack{z=(x-a_N)/b_N \text{ fixed}}]{x \rightarrow \infty, N \rightarrow \infty} F\left(\frac{x - a_N}{b_N}\right)$$

equivalently ,  $\lim_{N \rightarrow \infty} Q_N(a_N + b_N z) = F(z)$



## Gnedenko's classical law of extremes

For the maximum statistics from i.i.d. variables only **three** possible **forms** of  $F(z)$  exist. The particular one depends only the large  $x$  tail of the source distribution  $p(x)$ .

**Fisher-Tippett-Gnedenko Theorem.** Let  $\{X_n\}$  be a sequence of i.i.d. random variables. If there exist norming constants  $a_n > 0$ ,  $b_n \in \mathbb{R}$  and some non-degenerate CDF  $G$  such that  $a_n^{-1} (M_n - b_n)$  converges in distribution to  $G$ , then  $G$  is of one of the following three CDFs:

$$\text{Gumbel: } G_0(x) = \exp(-e^{-x}), \quad x \in \mathbb{R},$$

$$\text{Fr\'echet: } G_{1,\alpha}(x) = \exp(-x^{-\alpha}), \quad x \geq 0, \alpha > 0,$$

$$\text{Reversed Weibull: } G_{2,\alpha}(x) = \exp(-(-x)^{-\alpha}), \quad x \leq 0, \alpha < 0.$$

$$G_\gamma(x) = \exp(-(1 + \gamma x)^{-1/\gamma}), \quad 1 + \gamma x > 0.$$

## Uncorrelated RVs. Fréchet distribution

For the maximum statistics from i.i.d. variables only **three** possible **forms** of  $F(z)$  exist. The particular one depends only the large  $x$  tail of the source distribution  $p(x)$ .

1)  $P(x)$  has a power-law tail  $p(x) \sim A x^{-(1+\alpha)}$

$$F_1(z) = \begin{cases} e^{-z^{-\alpha}} & \text{for } z \geq 0 \\ 0 & \text{for } z \leq 0 \end{cases} \quad \text{Fréchet distribution}$$

$$f_1(z) = F'_1(z) = \frac{\alpha}{z^{\alpha+1}} e^{-z^{-\alpha}}, \quad z \in [0, \infty)$$

TABLE 1. A list of distributions in the Fréchet domain

Distribution	$1 - F(x)$	Extreme value index
Pareto	$\sim Kx^{-\alpha}, K, \alpha > 0$	$\frac{1}{\alpha}$
$F(m, n)$	$\int_x^\infty \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \omega^{\frac{m}{2}-1} (1 + \frac{m}{n}\omega)^{-\frac{m+n}{2}} d\omega$ $x > 0; m, n > 0$	$\frac{2}{n}$
Fréchet	$\frac{1 - \exp(-x^{-\alpha})}{x > 0; \alpha > 0}$	$\frac{1}{\alpha}$
$T_n$	$\int_x^\infty \frac{2\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})} \left(1 + \frac{\omega^2}{n}\right)^{-\frac{n+1}{2}} d\omega$ $x > 0; m, n > 0$	$\frac{1}{n}$

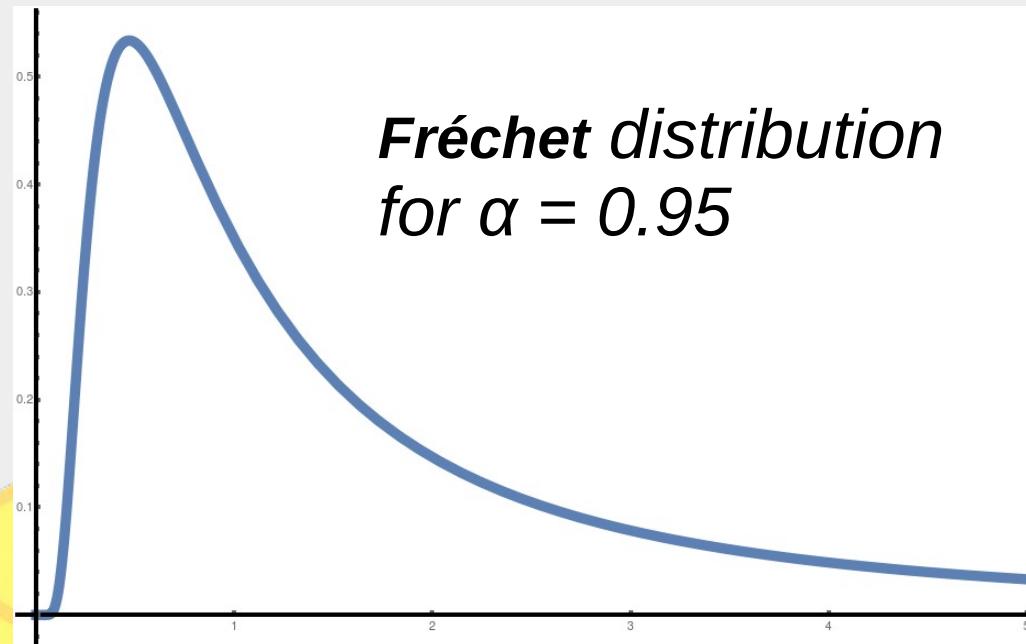
## Uncorrelated RVs. Fréchet distribution

For the maximum statistics from i.i.d. variables only **three** possible **forms** of  $F(z)$  exist. The particular one depends only the large  $x$  tail of the source distribution  $p(x)$ .

1)  $P(x)$  has a power-law tail  $p(x) \sim A x^{-(1+\alpha)}$

$$F_1(z) = \begin{cases} e^{-z^{-\alpha}} & \text{for } z \geq 0 \\ 0 & \text{for } z \leq 0 \end{cases}$$

$$f_1(z) = F'_1(z) = \frac{\alpha}{z^{\alpha+1}} e^{-z^{-\alpha}}, \quad z \in [0, \infty)$$



## Uncorrelated RVs. **Gumbel distribution**

2)  $p(x)$  has a tail decaying faster than a power law and unbounded

$$p(x) \sim e^{-x^\delta} \text{ with } \delta > 0$$

$$F_2(z) = e^{-e^{-z}}$$

### **Gumbel distribution**

$$f_2(z) = F'_2(z) = e^{-z - e^{-z}}, \quad z \in (-\infty, \infty)$$

TABLE 3. A list of distributions in the Gumbel domain

Distribution	$1 - F(x)$
Weibull	$\exp(-\lambda x^\tau), \quad x > 0; \lambda, \tau > 0$
Exponential	$\exp(-\lambda x), \quad x > 0; \lambda > 0$
Gamma	$\frac{\lambda^m}{\Gamma(m)} \int_x^\infty u^{m-1} \exp(-\lambda u) du, \quad x > 0; \alpha, m > 0$
Logistic	$1/(1 + \exp(x)), \quad x \in \mathbb{R}$
Normal	$\int_x^\infty \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) du, \quad x \in \mathbb{R}; \sigma > 0, \mu \in \mathbb{R}$
Log-normal	$\int_x^\infty \frac{1}{\sqrt{2\pi\sigma^2}u} \exp\left(-\frac{1}{2\sigma^2}(\log u - \mu)^2\right) du, \quad x > 0; \mu \in \mathbb{R}, \sigma > 0$

Q. Could you obtain  $F_2(z)$  for the exponential distribution  $p(x)=e^{-x}, x \geq 0$  ?

Reminder:  $F(z)$  is obtained from  $Q_N(x)$  which is

$$Q_N(x) = \left[ \int_{-\infty}^x dy \ p(y) \right]^N = \left[ 1 - \int_x^\infty dy \ p(y) \right]^N$$

$$Q_N(x) \xrightarrow[x \rightarrow \infty, N \rightarrow \infty]{z=(x-a_N)/b_N \text{ fixed}} F\left(\frac{x - a_N}{b_N}\right)$$

equivalently ,  $\lim_{N \rightarrow \infty} Q_N(a_N + b_N z) = F(z)$

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

Bernoulli law for Gaussian distr.  
i.i.d.

$$Q_N(x) = \left[ \int_{-\infty}^x dy p(y) \right]^N = \left[ 1 - \frac{1}{2} \operatorname{erfc} \left( \frac{x}{\sqrt{2}} \right) \right]^N = \exp \left[ N \ln \left( 1 - \frac{1}{2} \operatorname{erfc} \left( \frac{x}{\sqrt{2}} \right) \right) \right]$$

$$\operatorname{erfc}(z) = \frac{2}{\sqrt{\pi}} \int_z^{+\infty} e^{-y^2} dy$$

Expanding the ln:

$$Q_N(x) \approx \exp \left[ -\frac{N}{2} \operatorname{erfc} \left( \frac{x}{\sqrt{2}} \right) \right]$$

for  $x \gg 1$

$$\operatorname{erfc}(x) \approx \frac{e^{-x^2}}{\sqrt{\pi} x} \Rightarrow Q_N(x) \approx \exp \left( -\frac{N}{\sqrt{2\pi}} \frac{e^{-x^2/2}}{x} \right)$$

$$e^{-e^{-z}}$$

We want to get to universal form for  $z = \frac{x - aN}{bN}$   $\rightarrow \frac{aN}{bN}$ ?

$$Q_N(z) = Q_N(x = aN + bNz) = \exp \left[ -\frac{N}{\sqrt{2\pi}} \frac{\exp \left( -\frac{1}{2} (aN^2 + 2bNzAN + bN^2 z^2) \right)}{aN + bNz} \right]$$

$$bN \ll aN$$

$$Q_N(x = aN + bNz) = e^{-\frac{N}{\sqrt{2\pi} aN} \exp \left( -\frac{1}{2} (aN^2 + 2bNzAN) \right)}$$

$$\Rightarrow \begin{cases} \frac{N}{\sqrt{2\pi} aN} e^{-\frac{1}{2} aN^2} = 1 \\ bN = 1/aN \end{cases}$$

$$\begin{cases} \frac{N}{\sqrt{2\pi} aN} \approx \sqrt{2\ln N} - \frac{\ln(\ln N)}{2\sqrt{2\ln N}} \\ bN = 1/aN \end{cases}$$

$$\begin{aligned}
 & \text{Gumbel law for } e^{-x} = p(x) \\
 Q_N(x) &= \left[ \int_{-\infty}^x p(x) dx \right]^N = [1 - e^{-x}]^N = e^{N \ln[1 - e^{-x}]} = e^{-Ne^{-x}} = \\
 &= e^{-e^{-(x - \ln N)}} \\
 & z = x - \ln N \rightarrow a_N = \ln N \\
 & b_N = 1
 \end{aligned}$$

Remark: distribution with  $p(x) \sim C e^{-x^\delta}$ ,  $\delta > 0$

$$\begin{aligned}
 a_N &\approx (\ln N)^{1/\delta} \\
 b_N &\approx \frac{1}{\delta} (\ln N)^{1/\delta - 1} \rightarrow F_2(z) = e^{-e^{-z/\delta}}
 \end{aligned}$$

## Uncorrelated RVs. **Gumbel distribution**

2)  $p(x)$  has a tail decaying faster than a power law and unbounded

$$p(x) \sim e^{-x^\delta} \text{ with } \delta > 0$$

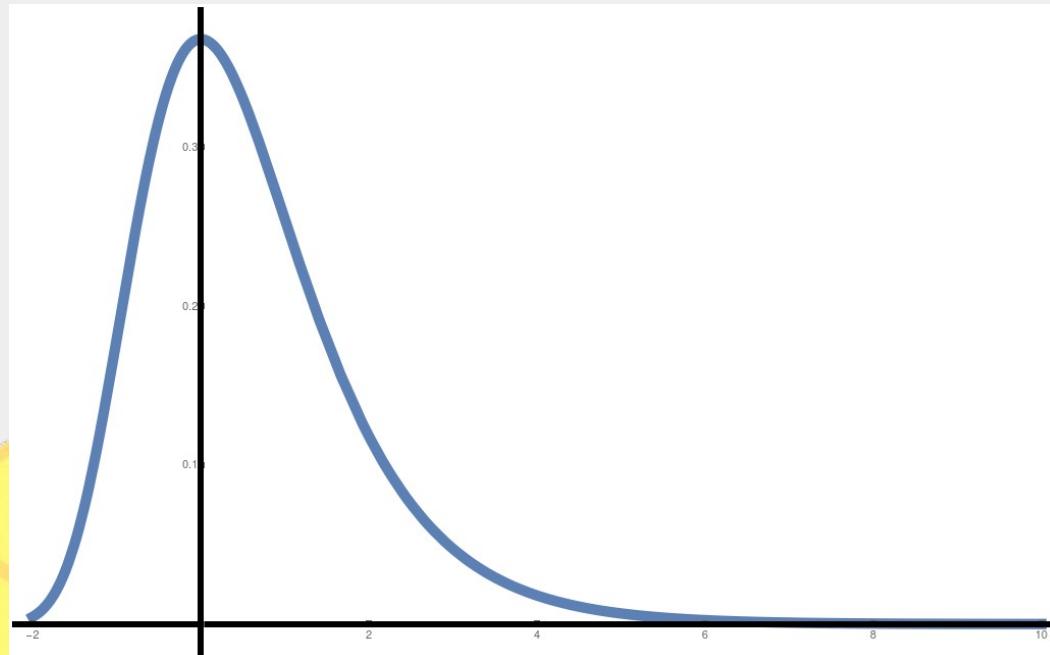
$$F_2(z) = e^{-e^{-z}}$$

### **Gumbel distribution**

$$f_2(z) = F'_2(z) = e^{-z - e^{-z}}, \quad z \in (-\infty, \infty)$$

TABLE 3. A list of distributions in the Gumbel domain

Distribution	$1 - F(x)$
Weibull	$\exp(-\lambda x^\tau), \quad x > 0; \lambda, \tau > 0$
Exponential	$\exp(-\lambda x), \quad x > 0; \lambda > 0$
Gamma	$\frac{\lambda^m}{\Gamma(m)} \int_x^\infty u^{m-1} \exp(-\lambda u) du, \quad x > 0; \alpha, m > 0$
Logistic	$1/(1 + \exp(x)), \quad x \in \mathbb{R}$
Normal	$\int_x^\infty \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) du, \quad x \in \mathbb{R}; \sigma > 0, \mu \in \mathbb{R}$
Log-normal	$\int_x^\infty \frac{1}{\sqrt{2\pi\sigma^2}u} \exp\left(-\frac{1}{2\sigma^2}(\log u - \mu)^2\right) du, \quad x > 0; \mu \in \mathbb{R}, \sigma > 0$



## *Uncorrelated RVs. Weibull distribution*

3)  $p(x)$  has a tail which is bounded, i.e.  $p(x) \xrightarrow{x \rightarrow a} (a - x)^{\beta - 1}$

$$F_3(z) = \begin{cases} e^{-(-z)^\beta} & \text{for } z \leq 0 \\ 1 & \text{for } z \geq 0 \end{cases}$$

**Weibull distribution**

$$f_3(z) = \beta(-z)^{\beta-1}e^{-(-z)^\beta}, \quad z \in (-\infty, 0]$$

TABLE 2. A list of distributions in the Reversed Weibull domain

Distribution	$1 - F(\omega(F) - \frac{1}{x})$	Extreme value index
Uniform	$\frac{1}{x}$ $x > 1$	-1
Beta( $p, q$ )	$\int_{1-\frac{1}{x}}^1 \frac{\Gamma(p+q)}{\Gamma(p)\Gamma(q)} u^{p-1} (1-u)^{q-1} du$ $x > 1; p, q > 0$	$-\frac{1}{q}$
Reversed Weibull	$1 - \exp(-x^{-\alpha})$ $x > 0; \alpha > 0$	$-\frac{1}{\alpha}$

# Skoltech Extreme value statistics

Skolkovo Institute of Science and Technology

*Uncorrelated RVs. Weibull distribution*

3)  $p(x)$  has a tail which is bounded, i.e.

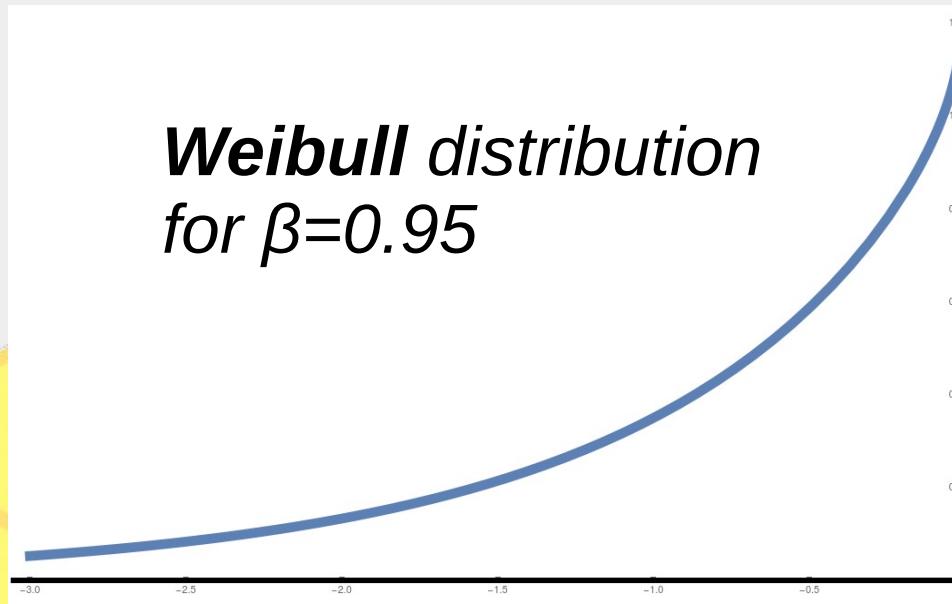
$$p(x) \xrightarrow{x \rightarrow a} (a - x)^{\beta - 1}$$

$$F_3(z) = \begin{cases} e^{-(-z)^\beta} & \text{for } z \leq 0 \\ 1 & \text{for } z \geq 0 \end{cases}$$

**Weibull distribution**

$$f_3(z) = \beta(-z)^{\beta-1}e^{-(-z)^\beta}, \quad z \in (-\infty, 0]$$

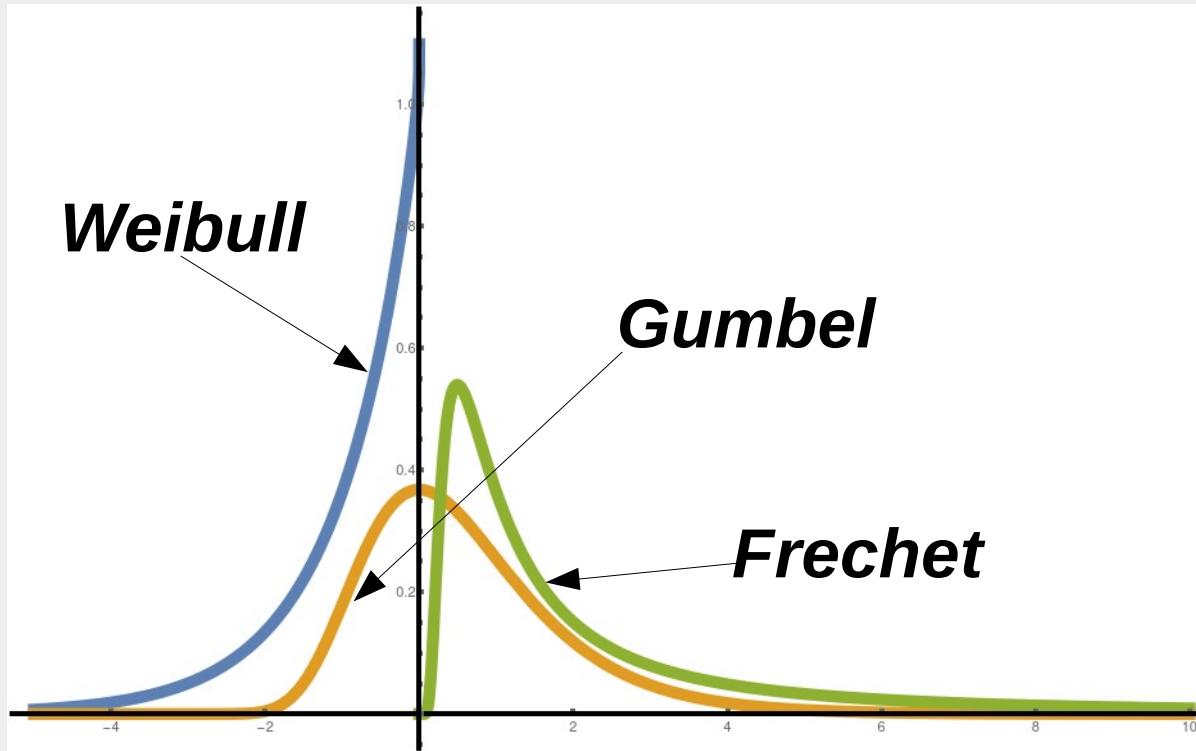
**Weibull distribution  
for  $\beta=0.95$**



# Skoltech Extreme value statistics

Skolkovo Institute of Science and Technology

*Uncorrelated RVs. All three universal distributions  
(Frechet, Gumbel, Weibull)*



Parent distribution $p(x)$	Scaling function $F(z)$	PDF of maximum $f(z)$	Nomenclature
$x^{-(1+\alpha)}; \alpha > 0$	$e^{-z^{-\alpha}} \theta(z)$	$\frac{\alpha}{z^{\alpha+1}} e^{-z^{-\alpha}}; z > 0$	<b>Fréchet</b>
$e^{-x^\delta}$	$e^{-e^{-z}}$	$e^{-z} e^{-e^{-z}}$	<b>Gumbel</b>
$(a - x)^{\beta-1}; \beta > 0$	$e^{-(-z)^\beta} \theta(-z) + \theta(z)$	$\beta(-z)^{\beta-1} e^{-(-z)^\beta}; z < 0$	<b>Weibull</b>

# Extreme value statistics

Example 1. Rogue (freak, killer) waves. Gumbel distr.



The Ninth Wave, I. Ayvazovsky, 1850



Rockall, 1943,  
by RAF photographer  
Fisher, James

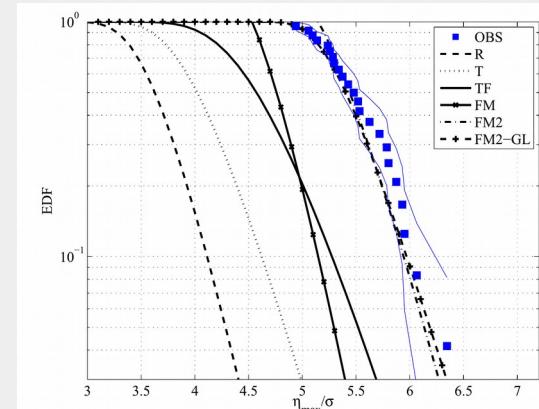
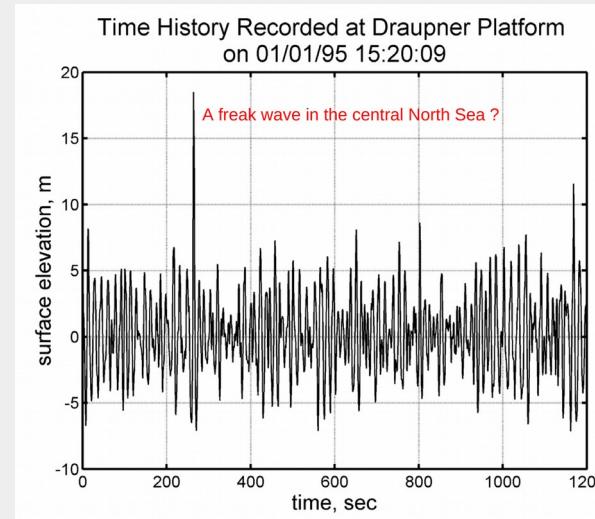


FIG. 9. Probability of exceedance (EDF) of dimensionless observed extreme crest heights  $\eta_{\max} = \max[\eta_i(t)]$  of the records  $B_i$  (OBS). The empirical EDF stability band is plotted as blue solid line. Reference distributions of extremes: Rayleigh (R); Tayfun (T); Tayfun-Fedele (TF); Fedele (FM); second-order FM (FM2); and asymptotic Gumbel limit of FM2 (FM2-GL). The reference duration is  $D = 1798$  s. Space-time extremes are computed over an area  $XY = 11.2 \times 11.2 \text{ m}^2 = 126 \text{ m}^2$ .

From P.H. Taylor, 1995

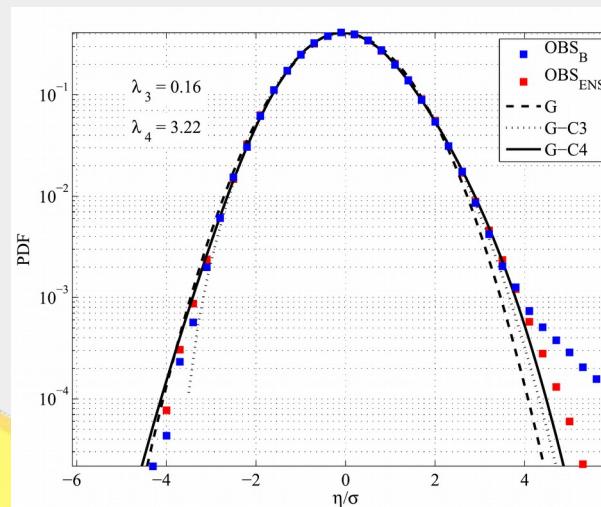


FIG. 7. PDF of dimensionless sea surface elevations of the records  $B_i$  ( $\text{OBS}_B$ ) and of the space-time ensemble ( $\text{OBS}_{\text{ENS}}$ ). The Gaussian distribution ( $G$ ), the third- ( $G\text{-C}3$ ), and fourth-order ( $G\text{-C}4$ ) nonlinear corrections are displayed for reference.

Benetazzo, A. et al.  
Observation of extreme  
sea waves  
in a space-time ensemble.  
J. Phys. Oceanogr. 45,  
2261–2275 (2015).



# Skoltech Extreme value statistics

Skolkovo Institute of Science and Technology

## Example 2. Financial markets. Frechet distribution

In financial markets the extreme price movements correspond to market corrections (ordinary periods) and stock market crashes during crisis periods.

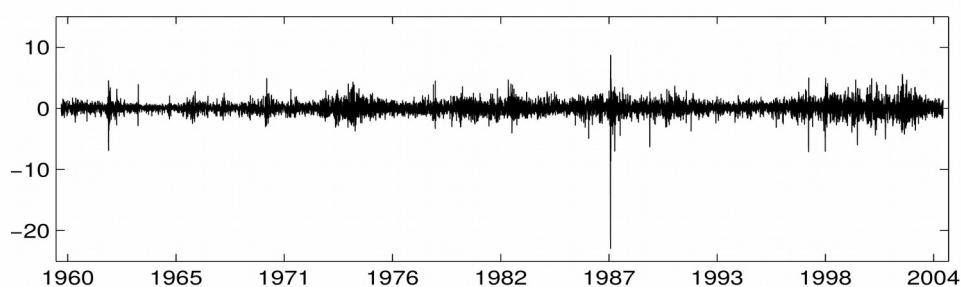


Fig. 5. Daily returns of the S&P500 index.

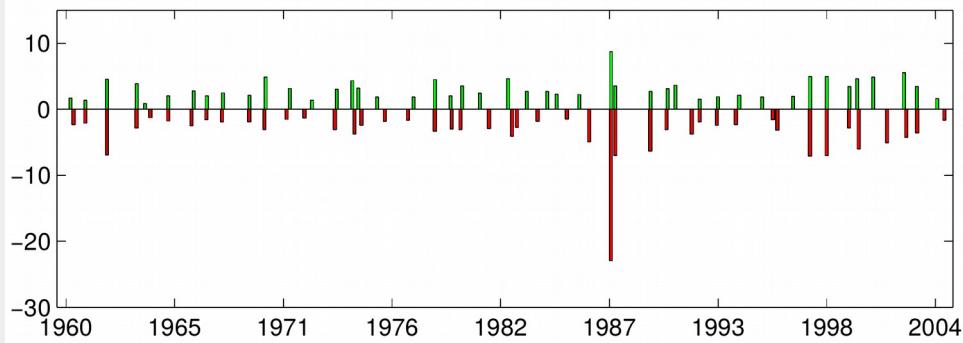


Fig. 6. Yearly minima and maxima of the daily returns of the S&P500.

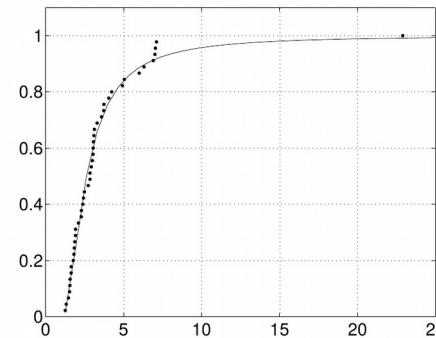
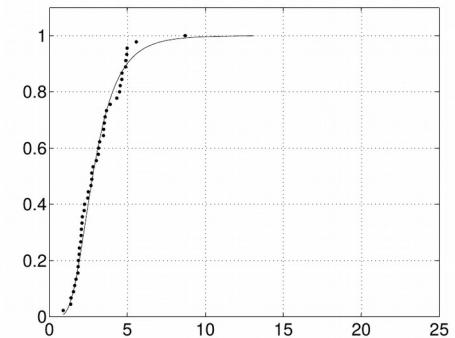
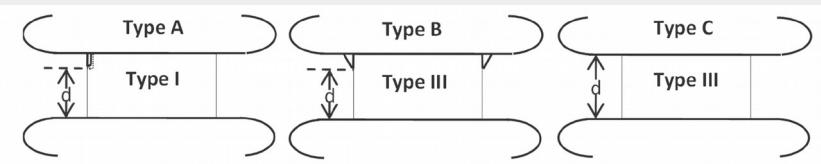


Fig. 7. Sample distribution (dots) of yearly minima (left panel) and maxima (right panel) and corresponding fitted GEV distribution for S&P500.



# Extreme value statistics

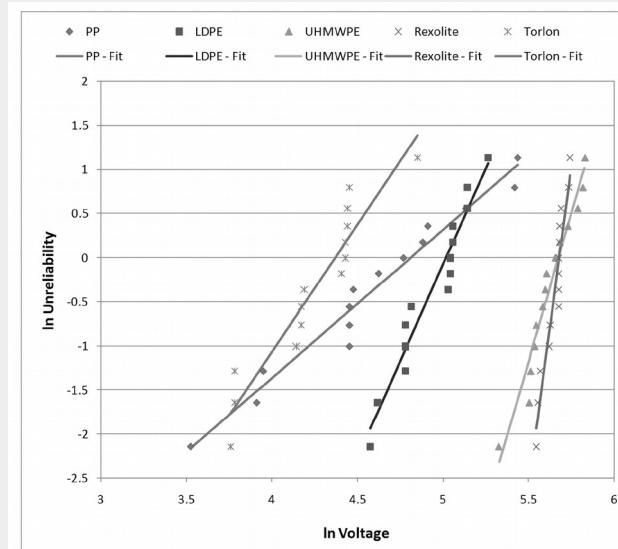
Example 3. Impulse-driven surface breakdown data.  
Weibull distribution



**Figure 1.** Sketch of the electrode/sample geometries tested and analysed (AI, BIII, and CIII).

Cumulative probability of failure

$$F(V) = 1 - \exp\left[-\left(\frac{V - \gamma}{\alpha}\right)^\beta\right]$$



**Figure 2.** Weibull plots and curve fits of peak applied voltage data for type I (recess) samples tested with high-voltage electrode type A (pin).

M.P. Wilson et al., 2011 IEEE Pulsed Power Conference,  
<https://doi.org/10.1109/PPC.2011.6191418>

# Order statistics

What is the statistics of the successive maxima and the gaps?

$$M_{1,N} = \max (x_1, x_2, \dots, x_N),$$

$$M_{2,N} = \text{second } \max (x_1, x_2, \dots, x_N),$$

⋮

$$M_{k,N} = k\text{-th } \max (x_1, x_2, \dots, x_N),$$

⋮

$$M_{N,N} = \min (x_1, x_2, \dots, x_N),$$

$$p_>(x) = \int_x^\infty p(y)dy$$

$$p_<(x) = \int_{-\infty}^x p(y)dy$$

$$Q_{k,N}(x) = \text{Prob}[M_{k,N} \leq x] = \sum_{m=0}^{k-1} \binom{N}{m} [p_>(x)]^m [p_<(x)]^{N-m}$$

$$Q_{k,N}(x) = \text{Prob}[M_{k,N} \leq x] \xrightarrow[x \rightarrow \infty, N \rightarrow \infty]{z = (x - a_N)/b_N \text{ fixed}} G_k \left( \frac{x - a_N}{b_N} \right)$$

$$G_k(z) = F_\mu(z) \sum_{j=0}^{k-1} \frac{[-\ln F_\mu(z)]^j}{j!}$$

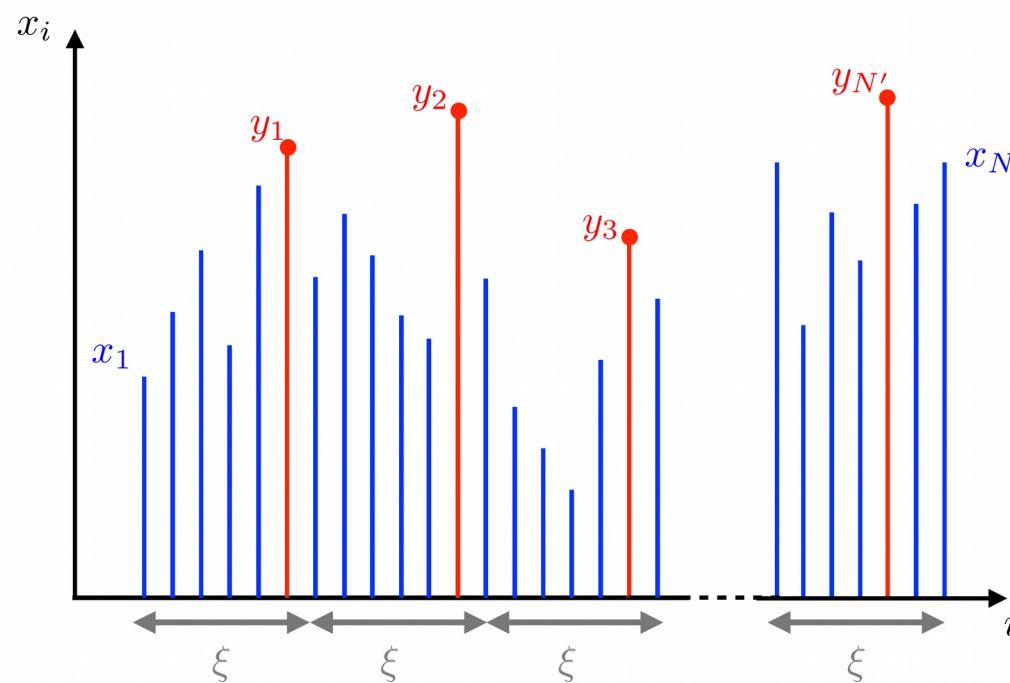
$F_\mu(z)$  is Frechet, Gumbel or Weibull

# Extreme value statistics

## Weakly correlated RVs

*Weak correlations* = finite correlation length. As a consequence the same universality as for i.i.d.

$$C_{i,j} = \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle \sim e^{-|i-j|/\xi}$$



$$M = \max[x_1, x_2, \dots, x_N] = \max[y_1, y_2, \dots, y_{N'}]$$

To decide the tail of  $p(y)$  one needs to solve a strongly correlated problem

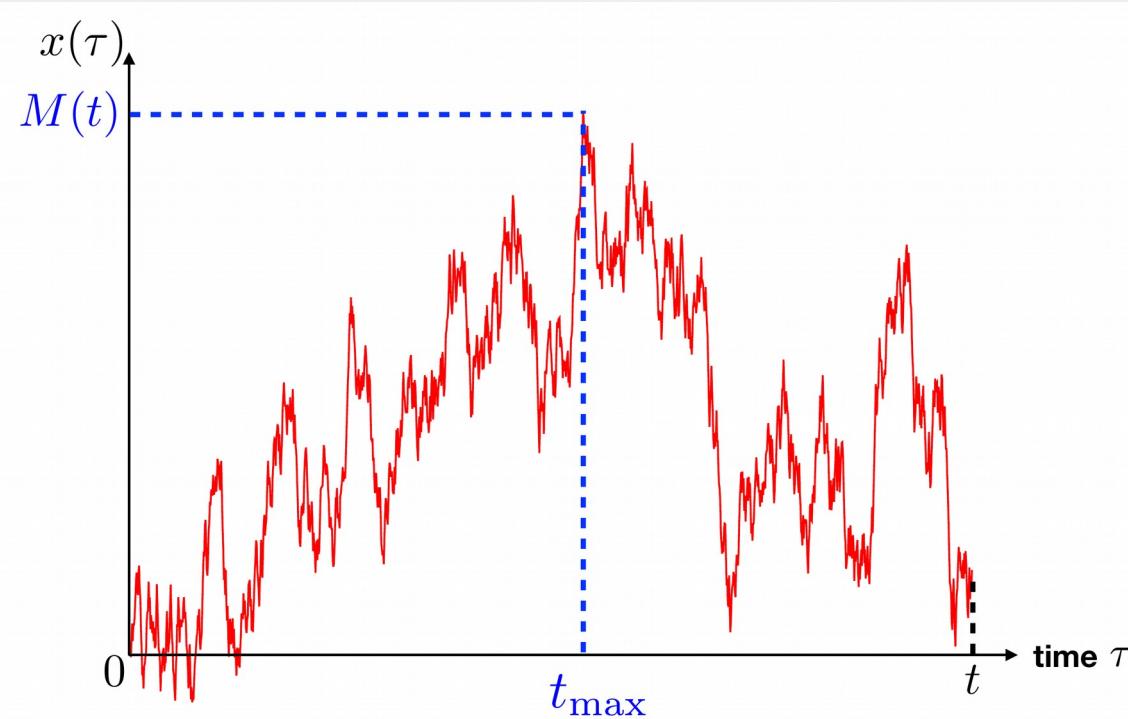
# Extreme value statistics

## Strongly correlated RVs

*Strong correlations*

No universality discovered yet. A field with many open questions.  
Some particular problems solved

Example: trajectory of Brownian motion



## 1D Brownian motion and RWs

$x_0 = 0$   
Langevin equation

$$\frac{dx}{dt} = \eta(t) \quad \text{noise}$$

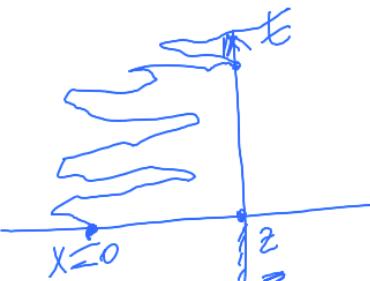
$$\left\{ \begin{array}{l} \langle \eta(t) \rangle = 0 \\ \langle \eta(t)\eta(t') \rangle = 2D\delta(t-t') \end{array} \right.$$

$$M(t) = \max_{0 \leq s \leq t} [x(s)]$$

$$x(t) = \int_0^t \eta(s) ds$$

white noise

$$\langle x(t) \rangle = 0, \quad \langle x(t)x(t') \rangle = 2D \min(t, t')$$



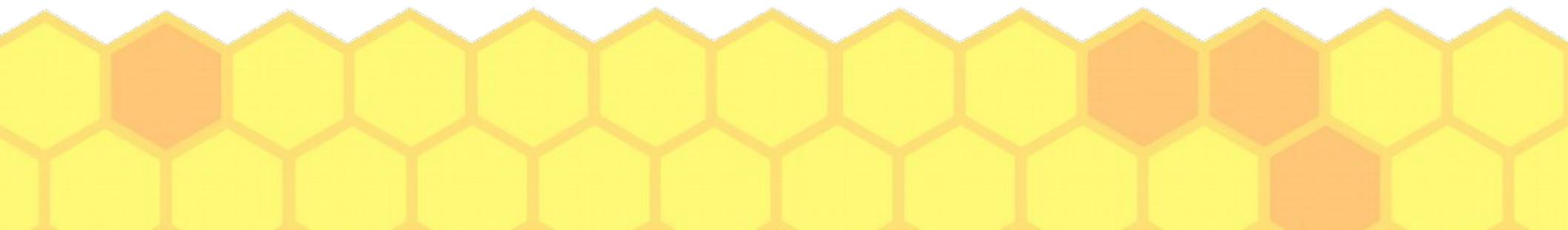
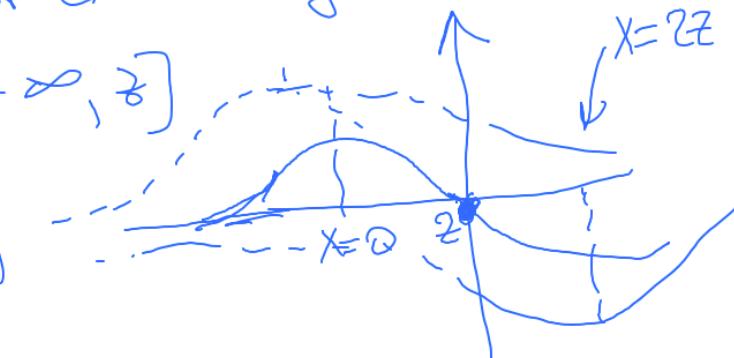
CDF of the maximum

$$Q(z,t) = \text{Prob}[M(t) \leq z] = \text{Prob}[x(s) \leq z, 0 \leq s \leq t]$$

$Q(z,0) = 0$   
 $P(x,t|z)$  is PDF that particle is at  $(x,t)$  not exceeding  $z$

$$\frac{\partial P}{\partial t} = D \frac{\partial^2 P}{\partial x^2} \quad x \in (-\infty, z]$$

$$P(x,t|z) = \frac{1}{\sqrt{4\pi Dt}} \left[ e^{-\frac{x^2}{4Dt}} - e^{-\frac{(x-z)^2}{4Dt}} \right]$$



$$Q(z,t) = \int_{-\infty}^z dx P(x,t|0) = \operatorname{erf}\left(\frac{z}{\sqrt{4Dt}}\right)$$

Heaviside function  $\operatorname{erf} = \frac{2}{\sqrt{\pi}} \int_{-\infty}^x e^{-y^2} dy$

$$P_M(z,t) = \frac{\partial Q(z,t)}{\partial z} = \Theta(z) \frac{1}{\sqrt{\pi Dt}} \exp\left(-\frac{z^2}{4Dt}\right)$$

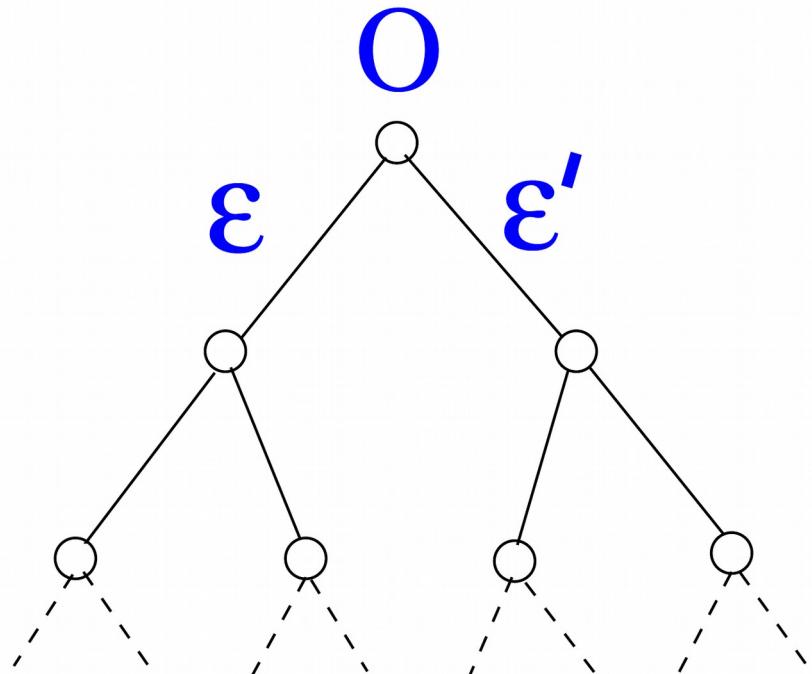
$$\langle z \rangle = M(t) = \int_0^\infty dz \frac{z}{\sqrt{\pi Dt}} e^{-\frac{z^2}{4Dt}} = \underbrace{\frac{2\sqrt{Dt}}{\sqrt{\pi}}}_{=} = M(t)$$

$$\langle \sigma x^2(t) \rangle = 2Dt \quad \text{for B.M}$$

# Extreme value statistics

Hierarchically correlated random variables

Disordered environment. Directed polymer on a Cayley tree



BST height. Maximum depth  
of the nodes in a tree  
with  $N$  occupied nodes

For BST:

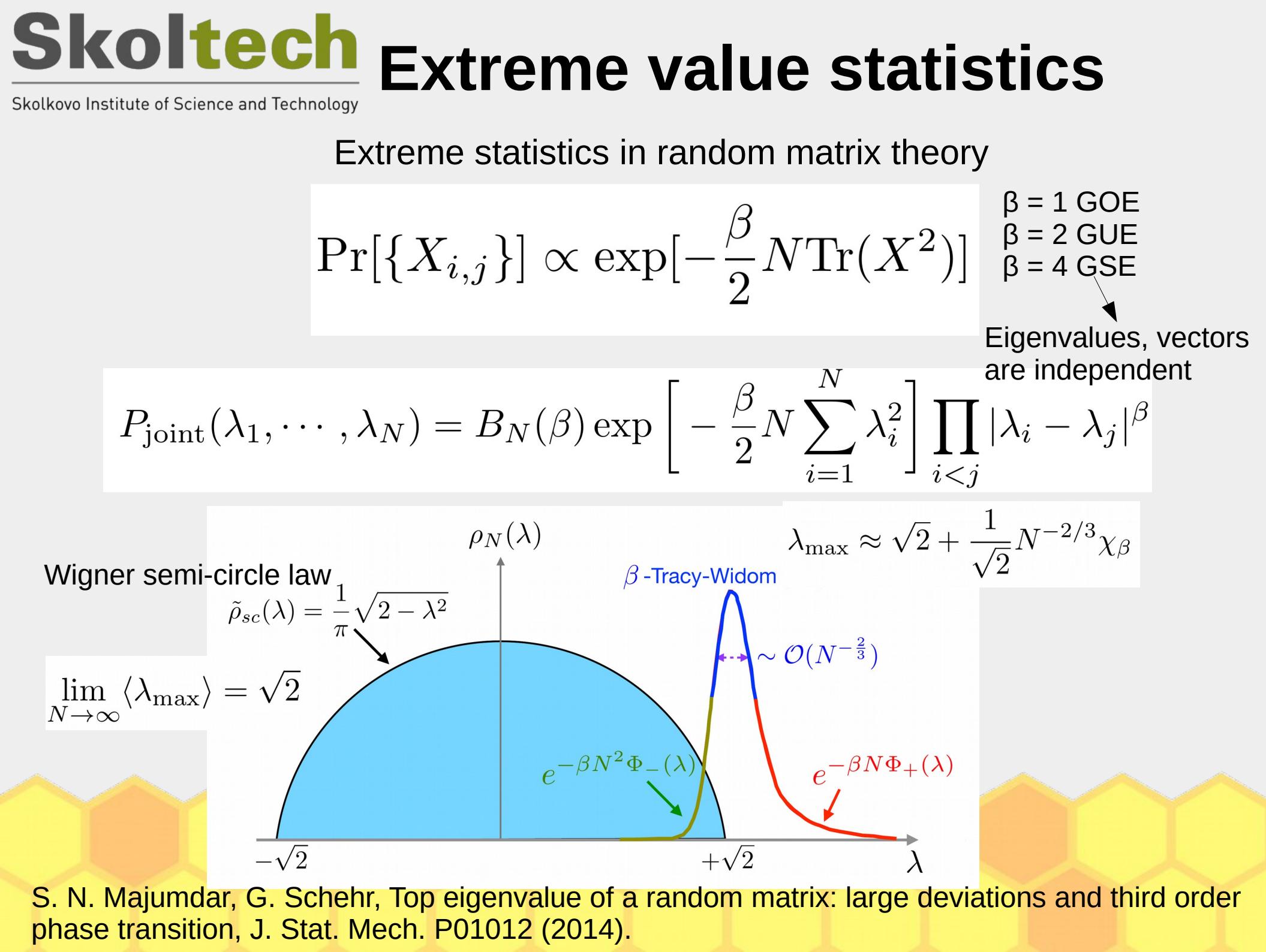
$$\langle E_{\min} \rangle \approx v(\lambda^*) n + \frac{3}{2\lambda^*} \ln n$$

$$v(\lambda) = -\frac{1}{\lambda} \ln \left[ 2 \int \rho(\epsilon) e^{-\lambda \epsilon} d\epsilon \right]$$

Equivalent to binary search tree (BST!!!)

If the data are stored on a tree of size  $N$  the optimal  
search time scaling is  $O(\ln N)$  not linear in  $N$

$$\langle H_N \rangle \approx \frac{1}{v(\lambda^*)} \ln N - \frac{3}{2\lambda^* v(\lambda^*)} \ln \ln N$$



# Extreme value statistics

## Other questions

- 1) Density of near-extremes
- 2) The time at which the max/min is reached
- 3) stats of records

Review of the topic with an emphasis of correlated RVs

**Extreme value statistics of correlated random variables: A pedagogical review,**

Satya N.Majumdar, Arnab Pal, Gr  gory Schehr, Physics Reports,

840, 2020, 1-32 or free version at <https://arxiv.org/abs/1910.10667>

(And references therein!)