

Entrega: curso de datos extremos

Laura Montaldo, CI: 3.512.962-7

2024-02-16



UNIVERSIDAD
DE LA REPÚBLICA
URUGUAY



Índice

Resumen	3
Motivación y objetivo del estudio	4
Marco Teórico	6
Teoría asintótica clásica y las distribuciones extremales y sus dominios de atracción	6
Definición 1: Las distribuciones extremales	7
Definición 2: Distribución extremal asintótica	12
Definición 3: Supremo esencial de una variable aleatoria o distribución	12
Definición 4: Distribución max-estables	14
Definición 4: Dominio de atracción maximal	17
Corolario 2 :	18
Definición 5: GEV	19
Referencias bibliográficas	20

Resumen

Your abstract goes here.

Motivación y objetivo del estudio

Los índices de *S&P* son una familia de índices de renta variable¹ diseñados para medir el rendimiento del mercado de acciones en Estados Unidos que cotizan en bolsas estadounidenses. Ésta familia de índices está compuesta por una amplia variedad de índices basados en tamaño, sector y estilo. Los índices están ponderados por el criterio *float-adjusted market capitalization* (FMC). Además, se disponen de índices ponderados de manera equitativa y con límite de capitalización de mercado, como es el caso del *S&P 500*. En este sentido, el *S&P500* entraría en el conjunto de índices ponderados por capitalización bursátil ajustada a la flotación (ver [S&P Dow Jones Indices](#)). El mismo mide el rendimiento del segmento de gran capitalización del mercado estadounidense. Es considerado como un indicador representativo del mercado de renta variable de los Estados Unidos, y está compuesto por 500 empresas constituyentes.

Se busca crear un indicador de una posible crisis bursátil. Como variable de referencia de toma la relación de precios al cierre de ayer sobre la de hoy

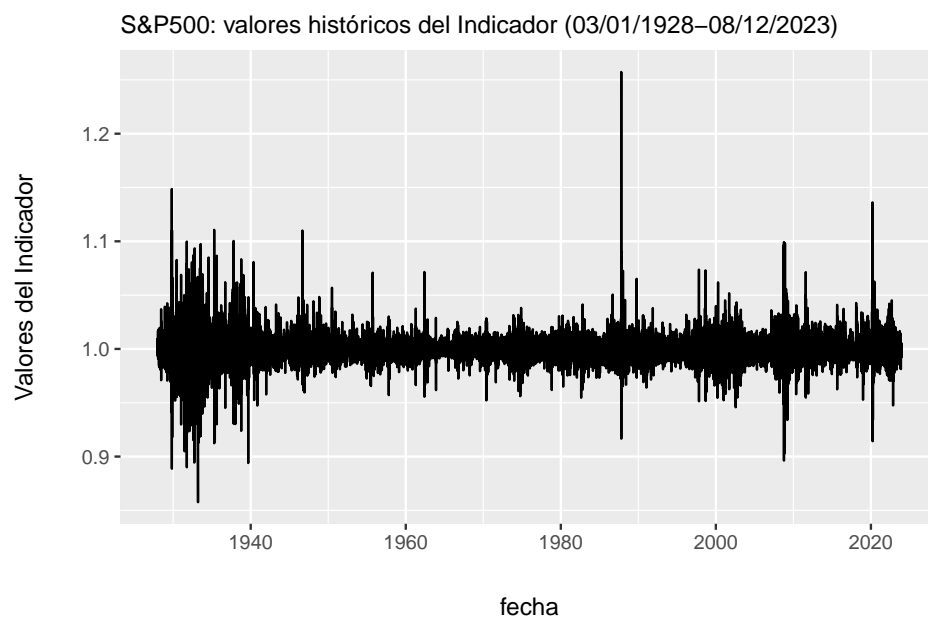
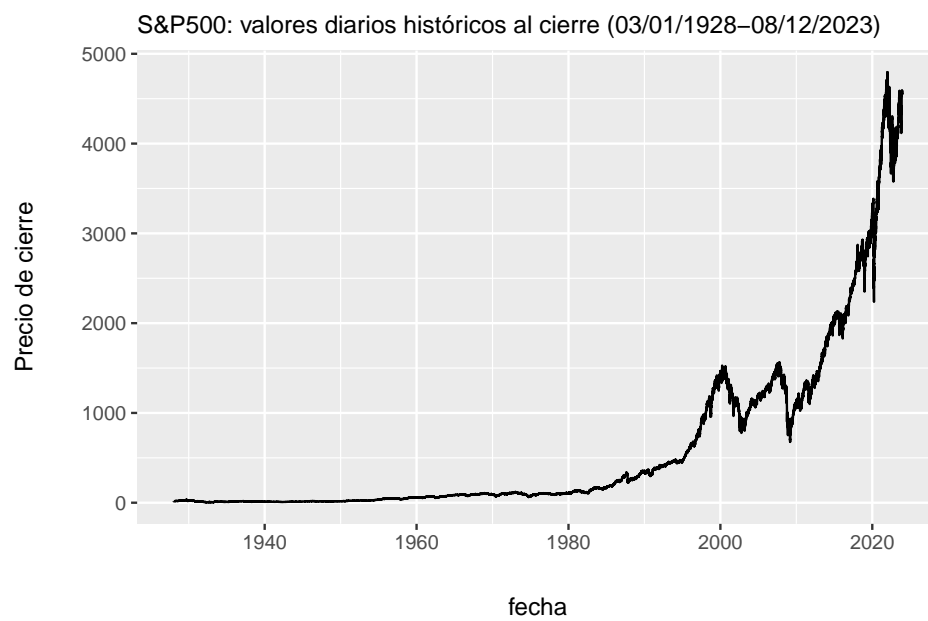
$$Indicador_t = \frac{Precio_{t-1}}{Precio_t}, \quad \text{para } t = 1, \dots, T \quad (1)$$

Interpretación del Indicador:

- Si el $Indicador_t \leq 1$, el precio de cierre de hoy es mayor o igual que el de ayer, lo cual podría ser considerado una señal positiva.
- Si el $Indicador_t > 1$, el precio de cierre de hoy es menor que el de ayer, lo cual podría considerarse una señal de alerta.

En la siguiente figura @ref(fig:plot1) se muestra la evolución histórica desde la fecha 03/01/1928 hasta 08/12/2023 del precio al cierre del día del indicador S&P 500.

¹En inglés se llaman equity indices



Marco Teórico

Teoría asintótica clásica y las distribuciones extremales y sus dominios de atracción

Siguiendo a Perera, Segura, y Crisci (2021) se dice que tenemos datos extremos cuando cada dato corresponde al máximo o mínimo de varios registros. Son un caso particular de evento raro o gran desviación respecto a la media.

Asumiremos que nuestros datos son *iid* (independientes e idénticamente distribuidos, son dos suposiciones juntas). Esta doble suposición suele no ser realista en aplicaciones concretas (ninguna de sus dos componentes, incluso) pero para comenzar a entender la teoría clásica, la utilizaremos por un tiempo.

Si tenemos datos X_1, \dots, X_n *iid* con distribución F , entonces $X_n^* = \max(X_1, \dots, X_n)$ tiene distribución F_n^* dada por $F_n^*(t) = F(t)_n$. Si conocemos la distribución F conoceríamos la distribución F_n^* , pero en algunos casos la lectura que queda registrada es la del dato máximo y no la de cada observación que dio lugar al mismo, por lo que a veces ni siquiera es viable estimar F . Pero aún en los casos en que F es conocida o estimable, si n es grande, la fórmula de F_n^* puede resultar prácticamente inmanejable. En una línea de trabajo similar a la que aporta el Teorema Central del Límite en la estadística de valores medios, un teorema nos va a permitir aproximar F_n^* por distribuciones más sencillas. Este es el Teorema de Fischer-Tippet-Gnedenko (FTG, para abreviar) que presentaremos en breve.

Como X_1, \dots, X_n *iid*, definimos $Y_i = -X_i$ para todo valor de i , entonces Y_1, \dots, Y_n *iid* y además $\min(X_1, \dots, X_n) = -\max(Y_1, \dots, Y_n)$ la teoría asintótica de los mínimos de datos *iid* se reduce a la de los máximos, razón por la que nos concentramos aquí en estudiar el comportamiento asintótico de los máximos exclusivamente.

Definición 1: Las distribuciones extremales

Las distribuciones extremales son tres: la distribución de Gumbel; la distribución de Weibull; la distribución de Fréchet.

Distribución de Gumbel Se dice que una variable tiene distribución de Gumbel si su distribución es:

$$\Lambda(x) = \exp\{-e^{-x}\} \quad \text{para todo } x \text{ real}$$

Distribución de Weibull Se dice que una variable tiene distribución de Weibull de orden $\alpha > 0$ si su distribución es:

$$\Psi_{\alpha}(x) = \begin{cases} \exp(-(-x)^{\alpha}) & \text{si } x < 0 \\ 1 & \text{en otro caso} \end{cases}$$

Distribución de Fréchet Se dice que una variable tiene distribución de Fréchet de orden $\alpha > 0$ si su distribución es:

$$\Phi_{\alpha}(x) = \begin{cases} \exp\{-x^{-\alpha}\} & \text{si } x > 0 \\ 0 & \text{en otro caso} \end{cases}$$

Teorema 1: Relaciones entre las versiones standard de las distribuciones extremales X tiene distribución $\Phi_\alpha(x)$ si y sólo si $(-1/X)$ tiene distribución $\Psi_\alpha(x)$ si y sólo si $\log(X^\alpha)$ tiene distribución Λ .

Teorema 2: Algunos datos de las distribuciones extremales

Parte 1 Si X tiene distribución $\Lambda^{(\mu, \beta)}$ entonces tiene:

- a) Valor esperado: $E(X) = \mu + \beta\gamma$, donde γ es la constante de Euler-Mascheroni, cuyo valor aproximado es 0.5772156649.
- b) Moda: μ
- c) Mediana: $\mu - \beta \log(\log 2) \approx \mu - 0.36651\beta$.
- d) Desviación estándar: $\beta\pi\sqrt{6} \approx 1.2825\beta$.
- e) Si $X^+ = \max(X, 0)$, entonces $E(X + k)$ es finito para todo valor de k natural.
- f) Para simular computacionalmente X , se puede tomar U uniforme en $(0, 1)$ y hacer $X = \mu - \beta \log(-\log U)$.

Parte 2 Si X tiene distribución $\Psi_\alpha^{(\mu, \beta)}$ entonces tiene:

- a) Valor esperado: $E(X) = \mu + \beta\Gamma(1 + 1/\alpha)$.
- b) Moda: μ si $\alpha \leq 1$ y $\mu - \beta\{(\alpha - 1)/\alpha\}^{(1/\alpha)}$ si $\alpha > 1$.
- c) Mediana: $\mu - \beta \log(2)^{(1/\alpha)}$.
- d) Desviación estándar: $\beta\{\Gamma(1 + 2/\alpha) - \Gamma(1 + 1/\alpha)^2\}^{1/2}$.

Parte 2 Si X tiene una distribución $\Phi_{\alpha}^{(\mu, \beta)}$ entonces se tiene:

- a) Valor esperado: $E(X) = \mu + \beta\Gamma(1 - 1/\alpha)$ si $\alpha > 1$, ∞ en caso contrario.
- b) Moda: $\mu + \beta\Gamma(1 - 1/\alpha)$ si $\alpha > 1$.
- c) Mediana: $\mu + \beta \log(2)^{(-1/\alpha)}$.
- d) Desviación estándar: $\beta|\Gamma(1 - 2/\alpha) - \Gamma(1 - 1/\alpha)^2|$ si $\alpha > 2$, ∞ si $1 < \alpha \leq 2$.

Teorema 3: Fischer-Tippet-Gnedenko (FTG) Si X_1, \dots, X_n *iid* con distribución F “continua”, llamamos F_n^* a la distribución de $\max(X_1, \dots, X_n)$ y n es grande, entonces existen μ real y $\beta > 0$ tales que alguna de las siguientes tres afirmaciones es correcta:

- 1) F_n^* se puede aproximar por la distribución de $\mu + \beta Y$ con Y variable con distribución Λ .
- 2) Existe $\alpha > 0$ tal que F_n^* se puede aproximar por la distribución de $\mu + \beta Y$ con Y variable con distribución Φ_α .
- 3) Existe $\alpha > 0$ tal que F_n^* se puede aproximar por la distribución de $\mu + \beta Y$ con Y variable con distribución Φ_α .

Lo anterior equivale a decir que la distribución del máximo de datos *continuos* e *iid*, si n es grande, puede aproximarse por una Gumbel, una Fréchet o una Weibull. Una aproximación será válida dependiendo de la distribución de F . En este sentido, cuando F sea normal entonces F_n^* se puede aproximar como una Gumbel. Cuando F sea uniforme, se puede aproximar F_n^* como una Weibull y cuando F sea Cauchy entonces F_n^* se puede aproximar por una Fréchet.

Más precisamente, cuál de las tres aproximaciones es la aplicable depende de la cola de F (los valores de $F(t)$ para valores grandes de t). En concreto, Weibull aparece cuando F es la distribución de una variable acotada por arriba (como la Uniforme), Gumbel para distribuciones de variables no acotadas por arriba pero con colas muy livianas (caso Exponencial y Normal) y Fréchet para colas pesadas (caso Cauchy)².

Como consecuencia del *FTG* cuando se tengan datos máximos, las distribuciones maximales podrían ser candidatas de uno de los ajustes si

- la cantidad de registros es lo suficientemente grande
- los registros son *iid* aunque con versiones más generales del *FTG* este supuesto puede no cumplirse

²Si bien la hipótesis de continuidad de F no es esencial, si F tiene la distribución Binomial o Poisson, por ejemplo, no se puede aplicar ninguna de las tres aproximaciones anteriores.

Como la mayoría de tests de ajustes suponen datos *iid*, se van a realizar dos tests de aleatoriedad³ a los datos:

- Runs up and down
- Spearman correlation of ranks

Se emplea la prueba de ajuste χ^2 que requiere seleccionar una partición más o menos arbitraria de la recta real de intervalos siendo importante que en cada intervalo haya una cantidad lo suficientemente importante de datos de la muestra. En este sentido, se pueden tomar como extremos de los intervalos los quintiles empíricos de la muestra. Cabe mencionar que este test requiere estimar parámetros por el método de Máxima Verosimilitud Categórica.

Cabe mencionar que para este estudio la distribución de la variable a incorporar en este estudio no tiene que ser degenerada, es decir $H(t) = 0$ ó $H(t) = 1$.

³En inglés se expresa como *randomness*

Definición 2: Distribución extremal asintótica

Si X_1, \dots, X_n es *iid* con distribución F diremos que H no-degenerada es la Distribución Extremal Asintótica (DEA) de F^4 , si existen dos sucesiones de números reales, d_n y $c_n > 0$, tales que la distribución de

$$\frac{\max(X_1, \dots, X_n) - d_n}{c_n} \quad (2)$$

tiende a H cuando n tiende a infinito.

Definición 3: Supremo esencial de una variable aleatoria o distribución

Si X tiene distribución F , se llama supremo esencial de X , denotado como M_X o, indistintamente, supremo esencial de F , denotado M_F a

$$M_X = M_F = \sup\{t/F(t) < 1\} \quad (3)$$

Observación:

- Si F es $U(a, b)$, $M_F = b$
- Si F es $Bin(m, p)$, $M_F = m$
- Si F es Normal, Exponencial, Cauchy o Poisson, M_F es infinito.

Teorema 4 Si X_1, \dots, X_n es *iid* con distribución F cualquiera, entonces, para n tendiendo a infinito,

$$X_n^* = M_F = \max(X_1, \dots, X_n) \text{ tiende a } M_F \quad (4)$$

Observación:

El resultado anterior vale incluso si M_F es infinito, pero si M_F es finito, como $X_n^* - M_F$ tiende a cero, por analogía con el Teorema Central del Límite para

⁴Lo que equivale a decir que F tiene DEA H .

promedios, buscaríamos una sucesión $c_n > 0$ y que tienda a cero de modo tal que $(X^*_n - M_F)/c_n$ tienda a una distribución no-degenerada y de allí surge buscar la DEA.

Teorema 5 Si F es una distribución con M_F finito, y para X con distribución F se cumple que

$$P(X = M_F) > 0$$

entonces F NO admite DEA.

Observación:

Si F es $Bin(m, p)$, $M_F = m$. Si X tiene distribución F , entonces $P(X = M_F) = P(X = m) = p_m > 0$, así que la distribución $Bin(m, p)$ NO admite DEA, no se puede aproximar la distribución del máximo de una muestra *iid* de variables $Bin(m, p)$.

El Teorema anterior es un caso particular del próximo.

Teorema 6 Si F es una distribución con M_F finito o infinito que admite DEA, y X tiene distribución F , entonces el límite cuando t tiende a M_F por izquierda de $P(X > t)/P(X \geq t)$ debe ser 1.

Observación:

- Si F es una distribución de Poisson de parámetro $\lambda > 0$, M_F es infinito.
- Si k es un natural, entonces:

$$\begin{aligned} \frac{P(X > k)}{P(X \geq k)} &= \frac{P(X \geq k+1)}{P(X \geq k)} \\ &= 1 - \frac{P(X = k)}{P(X \geq k)} \approx 1 - \left(1 - \frac{\lambda}{k}\right) \end{aligned} \tag{5}$$

que tiende a 0 cuando k tiende a infinito, por lo cual F NO admite DEA, o sea que no se puede aproximar el máximo de una sucesión *iid* de variables de Poisson.

Observación:

El Teorema 6 brinda una condición NECESARIA pero NO SUFICIENTE para DEA. Un ejemplo de ello lo aportó Von Mises, mostrando que la distribución

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

cumple con la condición del Teorema 6 pero no admite DEA.

Definición 4: Distribución max-estables

Si dada una F distribución, X con distribución F , k natural arbitrario y X_1, \dots, X_k es *iid* con distribución F , existen reales a_k, b_k tales que $\max(X_1, \dots, X_k)$ tiene la misma distribución que $a_k X + b_k$, F se dice *max-estable*.

El Teorema FTG resulta de superponer los dos siguientes teoremas:

Teorema 7

- a) Si F admite *DEA* H , entonces H es max-estable.
- b) Si H es max-estable, es la *DEA* de sí misma.

Teorema 8 Una distribución es max-estable si y solo si es extremal⁵. El Teorema 7 es bastante intuitivo y análogo a los teoremas de Lévy sobre distribuciones estables en aproximaciones asintóticas de las distribuciones de sumas. Para el Teorema 8 haremos enseguida un ejercicio sencillo que nos ayudará a hacerlo creíble. Luego precisaremos, para terminar con esta parte, cómo son las distribuciones que tienen por *DEA* cada uno de los tres tipos de distribuciones extremales. Para eso precisamos recordar algunas definiciones, como la siguiente.

Obsrvación:

Si F y G son dos distribuciones, tienen colas equivalentes si $M_F = M_G$ y cuando t tiende a M_F por izquierda, $(1 - F(t))/(1 - G(t))$ tiende a un

⁵O sea Gumbel, Weibull o Fréchet

valor $c > 0$. Recordando ahora cómo se calcula la distribución del máximo de dos variables independientes, es muy sencillo calcular la distribución del $\max\{X, Y\}$, cuando X e Y son independientes y cada una de ellas es una distribución extremal.

Se tiene el siguiente resultado:

X	Y	$\max(X, Y)$
Weibull	Weibull	Weibull
Weibull	Gumbel	Cola equivalente Gumbel
Weibull	Fréchet	Fréchet
Gumbel	Weibull	Cola equivalente Gumbel
Gumbel	Gumbel	Gumbel
Gumbel	Fréchet	Cola equivalente Fréchet
Fréchet	Weibull	Fréchet
Fréchet	Gumbel	Cola equivalente Fréchet
Fréchet	Fréchet	Fréchet

■ Las extremales son max-estables: tomar máximos de dos del mismo tipo queda en el mismo tipo.

■ Gumbel es más pesada que Weibull. En la cola, que es lo que cuenta para máximos, prima Gumbel.

■ Fréchet es más pesada que Gumbel y mucho más pesada que Weibull.

Además, de la tabla se deduce que

Teorema 9 Si X_1, \dots, X_n independientes y cada X_i tiene uno de los tres tipos de distribución extremal, entonces la distribución del $\max(X_1, \dots, X_n)$ es:

- Cola equivalente a Fréchet, si alguna de las variables es Fréchet y alguna otra es Gumbel.
- Fréchet, si alguna es Fréchet y ninguna es Gumbel.

- c) Cola equivalente Gumbel ninguna es Fréchet pero algunas son Gumbel y otras Weibull.
- d) Gumbel si todas son Gumbel.
- e) Weibull si todas son Weibull.

Observación:

Si F es una distribución, se dice que tiene *cola de variación regular de orden* $-\alpha$, para $\alpha \geq 0$, si para todo $t > 0$, $(1 - F(tx))/(1 - F(x))$ tiende a $t^{-\alpha}$ si $x \rightarrow \infty$. Para abreviar se dirá que F es $R_{-\alpha}$. Por ejemplo, para $\alpha = 3$, un caso de una tal F es $F(u) = 1 - 1/u^3$.

Por otra parte se dice que L es una *función de variación lenta* si, para todo $t > 0$, $L(tx)/L(x)$ tiende a 1 cuando $x \rightarrow \infty$. Por ejemplo, $L(u) = \log(u)$.

Definición 4: Dominio de atracción maximal

Si H es una distribución extremal (Gumbel, Weibull o Fréchet) su Dominio de Atracción Maximal ($DAM(H)$) está constituido por todas las distribuciones F que tienen $DEA H$.

Teorema 9: DAM de la Fréchet F pertenece a la DAM de Φ_α si y sólo si $1 - F(x) = x - \alpha L(x)$ para alguna L de variación lenta, lo cual es equivalente a decir que F es $R_{-\alpha}$.

Corolario 1: DAM de la Fréchet Si F es una distribución con densidad f que cumple que $xf(x)/(1 - F(x))$ tiende a α cuando $x \rightarrow \infty$ se dice que F cumple la Condición de Von Mises I. En tal caso, F pertenece a la DAM de Φ_α y mas aún, la DAM de Φ_α son todas las distribuciones que tienen cola equivalente a alguna distribución que cumpla la Condición de Von Mises I. Del DAM Fréchet y Teorema 1, surge lo siguiente.

Teorema 10: DAM de la Weibull

a) F pertenece a la DAM de Ψ_α si y solo si M_F es finito y además

$$1 - F(M_F - 1/x) = x^{-\alpha} L(x)$$

para alguna L de variación lenta, es decir que pertenece a $R_{-\alpha}$. Observar que con el cambio de variable $u = M_F - 1/x$, resulta que $1 - F(u) = (-MF - u)^\alpha L(1/(M_F - u))$ para alguna L de variación lenta, para $u < M_F$. Además puede tomarse $d_n = M_F$ y $c_n = n - \alpha$.

b) Si F distribución con densidad f positiva en (a, M_F) para algun $a < M_F$ y $(M_F - x)f(x)/(1 - F(x))$ tiende a α cuando $x \rightarrow M_F$, se dice que F cumple la Condición de Von Mises II. En tal caso F pertenece a la DAM de Ψ_α y mas aún, la DAM de Ψ_α son todas las distribuciones que tienen cola equivalente a alguna distribución que cumpla la Condición de Von Mises II.

Teorema 11: DAM de la Gumbel Una distribución F se dice una Función de Von Mises con función auxiliar h si existe $a < M_F$ (M_F puede ser finito o infinito) tal que para algún $c > 0$ se tiene

$$1 - F(x) = c \exp^{-\int_a^x \frac{1}{h(t)} dt},$$

con h positiva, con densidad h' y $h'(x)$ tendiendo a 0 para $x \rightarrow M_F$. Se tiene entonces que la DAM de Λ son todas las distribuciones que tienen cola equivalente a alguna distribución que sea una Función de Von Mises. Básicamente, se trata de colas más livianas que cualquier expresión del tipo $1/x^k$, más aún, con decaimiento *del tipo exponencial*, en el sentido preciso siguiente: si como en el Teorema 11

$1 - F(x) = c \exp^{-\int_a^x \frac{1}{h(t)} dt}$, entonces se tiene $1 - F(x) = c \exp^{-(x-a)/h(x)}$, donde la función auxiliar h es no-decreciente y con asíntota horizontal.

Además, d_n y c_n suelen involucrar expresiones logarítmicas. Más concretamente, $d_n = F^{-1}(1 - 1/n)$, $c_n = h(d_n)$, donde F^{-1} es la inversa generalizada (o función cuantil), definida por $F^{-1}(p) = \inf\{t/F(t) \geq p\}$, para $0 < p < 1$.

Corolario 2 :

Si F pertenece al DAM Gumbel, M_F es infinito, y se considera X con distribución F , entonces $E(X + k)$ es finito para todo k natural. Los resultados antes vistos nos permiten reconocer que distribuciones tienen DEA y si la tienen, cual es. Cierran el tema. Adicionalmente, permiten ver con mucha precisión que el quid de esta teoría es el comportamiento de las colas de las distribuciones, que Fréchet corresponde a las colas más pesadas, luego la Gumbel y finalmente Weibull. Para terminar el capítulo presentaremos la distribución de valores extremos generalizada⁶, que es una forma de compactar en una única fórmula las tres distribuciones extremas, debida a Jenkinson-Von Mises.

⁶GEV, por sus siglas en inglés.

Definición 5: GEV

Se define la distribución GEV de posición μ , escala β e índice ξ a

El caso $\xi=0$ corresponde a Gumbel, el caso $\xi < 0$ a Weibull y $\xi = -1/$, el caso $\xi > 0$ a Fréchet y $\xi = 1/$ En R existen rutinas para estimar con intervalos de confianza (por máxima verosimilitud, etc.) lo cual da formas de testear si una extremal es Gumbel, Weibull o Fréchet.

Referencias bibliográficas

Perera, Gonzalo, Angel Segura, y Carolina Crisci. 2021. *Curso de estadística de datos extremos, cap. 1 a cap. 5*.