

Análisis de Datos

Métodos Numéricos y
Simulación.

Segundo de Grado en Física.

Objetivo del tema

1) En ocasiones, cuando uno hace una medida de una magnitud X no espera obtener el mismo resultado aunque la medida se repita exactamente en las mismas condiciones.

Se dice entonces que X constituye una variable aleatoria

Aprenderemos a describir series de datos de variables aleatorias y a sacar conclusiones de ellas.

2) En otras ocasiones, se observa que los valores que toma una magnitud Y , aunque tienen cierta aleatoriedad, se ven influidos por los valores de otras magnitudes X, Z, T, \dots , y se dispone de un modelo $f(X, Z, T, \dots)$ para esa posible influencia que depende de una serie de parámetros.

Aprenderemos a comprobar si la influencia de X, Z, T, \dots es significativa y a encontrar los valores óptimos de dichos parámetros.

Importar datos a Matlab

El primer paso es importar los datos a Matlab.

- Copiando y pegando.
- Usando el *Import Data Wizard* (Menú File/Import Data)
- Mediante comandos en un script:
 - fid = fopen(NombreArchivo)***: Abre el archivo *NombreArchivo* y le asigna el identificador *fid*
 - fclose(fid)***: Cierra el archivo con identificador *fid*.
 - C=textscan(fid,format)***: lee del archivo con identificador *fid* mientras los contenidos se ajustan al formato especificado en *format*. En caso contrario, se detiene. Los datos leídos se almacenan en una matriz de celdas *C*.
 - tline=fgetl(fid)***: lee una línea de texto del archivo con identificador *fid* y la almacena en la cadena de caracteres *tline*.
 - otros**: *textread*, *csvread*, *fread*, *fscanf*, etc...

Ejemplo (importar datos a Matlab)

Importar los datos del radio y carga eléctrica de partículas de un medio granular en suspensión.

```
4 - fid=fopen('Serie-ejemplo.txt');
5 - TextoInicial = fgetl(fid); % Lee la primera línea del archivo
6 - linea1 = fgetl(fid);
7 - linea2 = fgetl(fid);
8 - linea3 = fgetl(fid); % Leen las tres líneas de texto antes de cada serie
9 - serie1 = textscan(fid,'%f %f'); % Lee mientras encuentre sólo los caracteres que aparecen entre corchetes
10 - linea1 = fgetl(fid);
11 - linea2 = fgetl(fid);
12 - linea3 = fgetl(fid); % Leen las tres líneas de texto antes de cada serie
13 - serie2 = textscan(fid,'%f %f'); % Lee mientras encuentre sólo los caracteres que aparecen entre corchetes
14 - linea1 = fgetl(fid);
15 - linea2 = fgetl(fid);
16 - linea3 = fgetl(fid); % Leen las tres líneas de texto antes de cada serie
17 - serie3 = textscan(fid,'%f %f'); % Lee mientras encuentre sólo los caracteres que aparecen entre corchetes
18 - linea1 = fgetl(fid);
19 - linea2 = fgetl(fid);
20 - linea3 = fgetl(fid); % Leen las tres líneas de texto antes de cada serie
21 - serie4 = textscan(fid,'%f %f'); % Lee mientras encuentre sólo los caracteres que aparecen entre corchetes
22 - linea1 = fgetl(fid);
23 - linea2 = fgetl(fid);
24 - linea3 = fgetl(fid); % Leen las tres líneas de texto antes de cada serie
25 - serie5 = textscan(fid,'%f %f'); % Lee mientras encuentre sólo los caracteres que aparecen entre corchetes
26 - fclose(fid);
27 - clear('TextoInicial','linea','linea1','linea2','linea3');
28 - r1=serie1{1};
29 - r2=serie2{1};
30 - r3=serie3{1};
31 - r4=serie4{1};
32 - r5=serie5{1};
33 - q1=serie1{2};
34 - q2=serie2{2};
35 - q3=serie3{2};
36 - q4=serie4{2};
37 - q5=serie5{2};
38 - clear('serie1','serie2','serie3','serie4','serie5','fid');
```

Algunas definiciones

- **Población:** el conjunto de todos los posible valores que puede tomar la magnitud objeto de nuestro estudio.

En nuestro caso, la población es el conjunto de las cargas y los radios de todas las partículas de polvo del bote.

- **Muestra:** el conjunto de datos que nosotros hemos medido.

- Para que lo que estudiemos en este tema sea válido, la muestra debe ser **aleatoria**.

En nuestro ejemplo, que la muestra sea aleatoria significa que todas las partículas del bote han tenido la misma probabilidad de participar en la muestra.

- Cuando una muestra no es aleatoria, se dice que está sesgada.

Descripción de una población

Cuando no es factible dar todos los valores de una población, ésta se describe dando:

- Su distribución acumulativa $F(x)$

$F(x_0)$: Fracción del total de medidas N en los que $x < x_0$.

- O bien su distribución diferencial $f(x)$

$f(x_0)\Delta x$ da la probabilidad de obtener un valor de x comprendido entre $x_0 - \Delta x$ u $x_0 + \Delta x$.

Propiedades:

$$\int_{-\infty}^{+\infty} f(x)dx = 1 \qquad F(x_0) = \int_{-\infty}^{x_0} f(x)dx \qquad f(x) = \frac{dF}{dx}$$

Descripción de una muestra

- No podemos determinar $F(x)$ y $f(x)$, porque eso significa medir para toda la población (en nuestro caso, para todos los granos del bote).
- Lo que podemos hacer es aproximar $F(x)$ y $f(x)$ a partir de las medidas experimentales.
- Voy a llamar $F_{\text{exp}}(x)$ a la distribución acumulativa y $f_{\text{exp}}(x)$ a la distribución diferencial obtenidas a partir de los datos experimentales
- $F_{\text{exp}}(x)$ y $f_{\text{exp}}(x)$ no se extienden entre $-\infty$ y $+\infty$, sino entre $x_{\text{min}} < x < x_{\text{max}}$, porque el número de valores de x que tenemos es finito)

Distribución acumulativa F_{exp}

- Supongamos que tengo una muestra x_1, x_2, \dots, x_N y quiero representar su distribución acumulativa.

1) Ordeno los x_j de menor a mayor: a la muestra así ordenada la llamo $x'_1, x'_2, \dots, x'_j, \dots, x'_N$

- Para ordenar una serie se usa la función **$\mathbf{B} = \text{sort}(\mathbf{A})$** .

A: vector desordenado.

B: vector ordenado en sentido ascendente.

2) Represento $1/N, 2/N, \dots, j/N, \dots, (N-1)/N, 1$ frente a $x'_1, x'_2, \dots, x'_j, \dots, x'_{N-1}, x'_N$

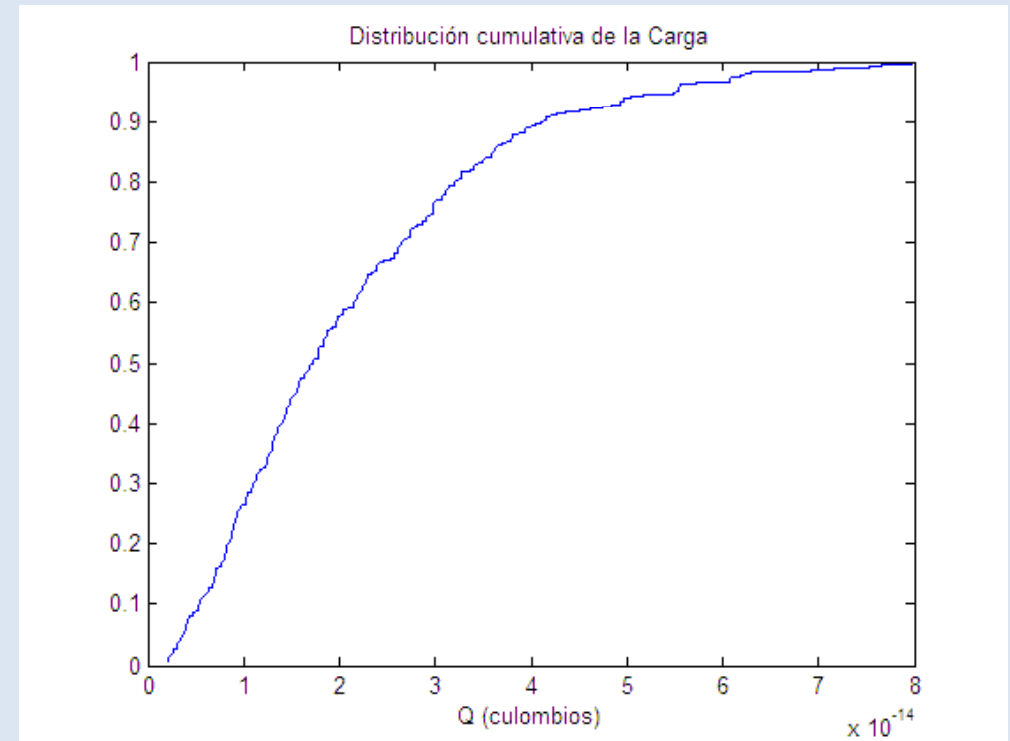
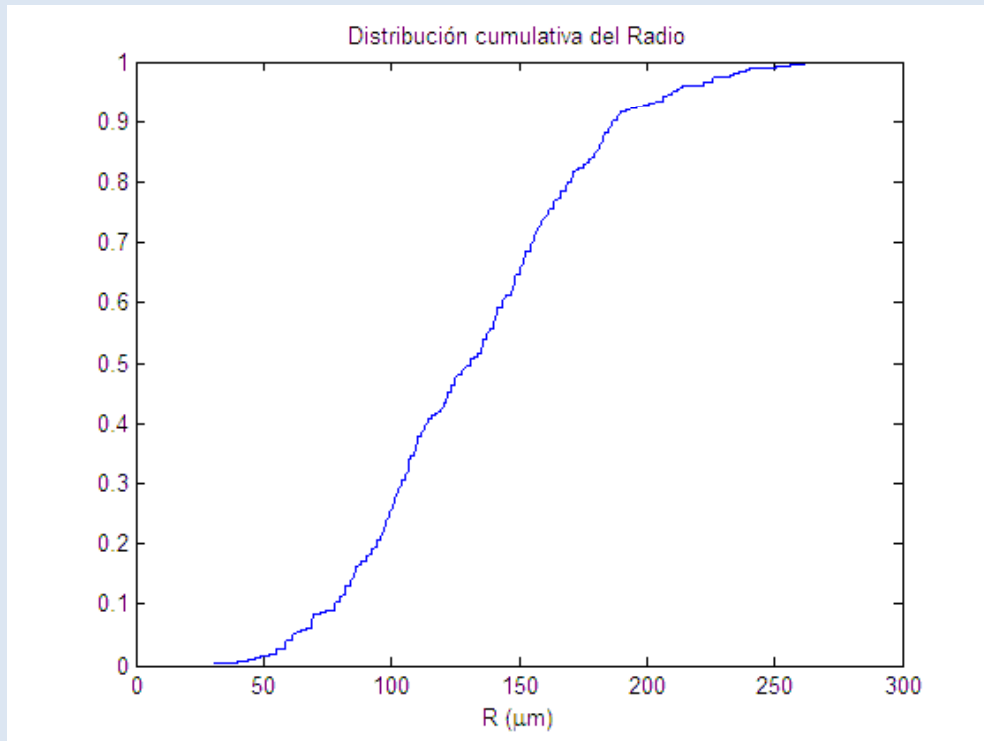
- Se usa la función **$\text{stairs}(\mathbf{X}, \mathbf{Y})$** .

Representa el vector Y frente al vector X haciendo cambiar Y sólo donde toma valores el vector X.

Ejemplo (distribución acumulativa)

```
Editor - F:\docencia\clases\materiales recurrentes\MNS\análisis\Cumulativas.m
File Edit Text Go Cell Tools Debug Desktop Window Help
[Icons] Stack: Base
[Icons] 1.0 1.1 % %
1 | % Dibuja la distribución acumulativa del radio y de la carga.
2 |
3 - r = [r1 ; r2 ; r3 ; r4 ; r5]; % Concateno todas las series del radio
4 - q = [q1 ; q2 ; q3 ; q4 ; q5]; % Y también de la carga
5 - R = sort(r); % Ordeno los radios de menor a mayor
6 - Q = sort(q);
7 - [N , M] = size(R);
8 - DistVal = 1:N;
9 - DistVal = DistVal/N;
10 - DistVal= DistVal'; % Transpongo para obtener un vector columna
11 - figure(1);
12 - stairs(R,DistVal);
13 - title('Distribución acumulativa del Radio');
14 - xlabel('R (\mum)');
15 - figure(2);
16 - stairs(Q,DistVal);
17 - title('Distribución acumulativa de la Carga');
18 - xlabel('Q (culombios)');
```

Ejemplo (distribución acumulativa)



Estas curvas representan $F_{\text{exp}}(R)$ y $F_{\text{exp}}(Q)$. Las funciones $F(R)$ y $F(Q)$ para nuestra población nos son desconocidas

Momentos de la distribución

Los momentos de una distribución nos ayudan a describir su forma numéricamente.

Son: la mediana, la media, la desviación estándar, la skewness y la curtosis.

Pueden calcularse tanto para la población como para la muestra que hemos medido.

Normalmente sólo podemos medir sus valores en una muestra: a partir de los valores de la muestra tenemos que estimar los de la población.

Media

Cálculo para la población:

$$\langle x \rangle = \int_{-\infty}^{+\infty} xf(x)dx$$

Normalmente no podemos calcularla de esta forma porque no conocemos $f(x)$. La estimamos a partir de una muestra de la forma:

$$\langle x \rangle = \frac{1}{N} \sum_{i=1}^N x_i$$

En Matlab, se usa la función:

$M = \text{mean}(x)$: M es la media de los elementos del vector x .

Mediana

Es el valor x_0 para el cual $F(x_0) = 1/2$, es decir, el 50% de los datos están por debajo de la mediana y 50% está por encima.

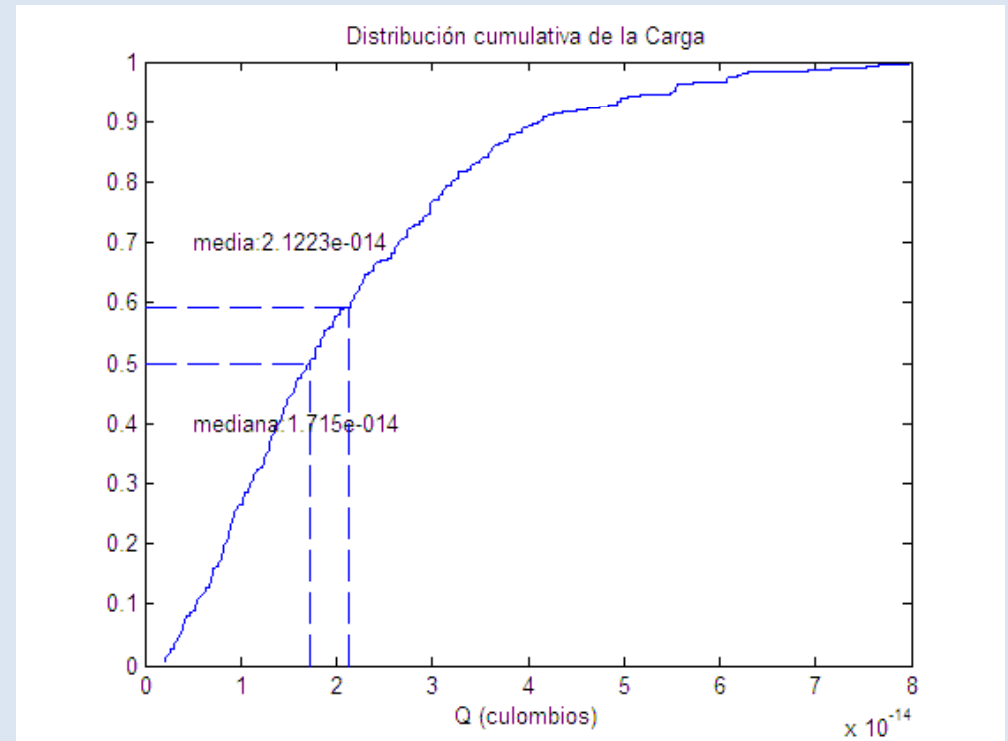
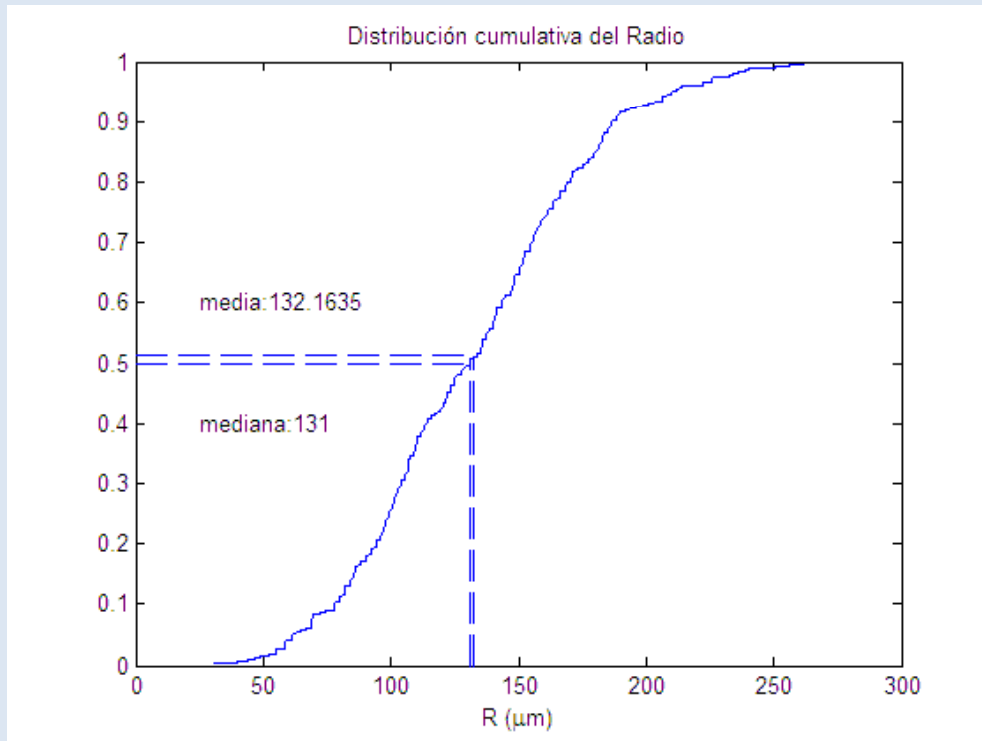
Para estimar su valor a partir de una muestra, buscamos $F_{\text{exp}}(x_0) = 1/2$.

en Matlab, se usa:

$M = \text{median}(x)$: M es la mediana de los valores en el vector X

Tiene la ventaja sobre la media de que resulta menos afectada por puntos erróneos (outliers)

Ejemplo (media y mediana)



Cuando una distribución no es simétrica, la mediana y la media no coinciden.

Varianza y desviación estándar

➤ Cálculo para la población:

$$\text{Var}_x = \int_{-\infty}^{+\infty} (x - \langle x \rangle)^2 f(x) dx$$

➤ Cálculo para una muestra:

$$\text{Var}_x = \frac{1}{N-1} \sum_{j=1}^N (x_j - \langle x \rangle)^2$$

A partir de la varianza, se define la desviación estándar como:

$$s(x), \sigma(x) = \sqrt{\text{Var}_x}$$

var(X, flag), ***std(X, flag)***: varianza y desviación estándar de los valores contenidos en el vector X . $flag = 0$ si los valores de X son una muestra de una población. $flag = 1$ si los valores de X corresponden a una población completa.

La desviación estándar da una medida de cuánto de agrupada está una distribución alrededor de su valor medio.

Histogramas de frecuencias (I)

Una forma cualitativa de dar $f_{\text{exp}}(x)$ es presentarla en la forma de un histograma de frecuencias. Para ello:

1) Defino una serie de M intervalos de anchura Δx :

$$\left(x_1^* - \frac{\Delta x}{2}, x_1^* + \frac{\Delta x}{2}\right), \left(x_2^* - \frac{\Delta x}{2}, x_2^* + \frac{\Delta x}{2}\right), \dots, \left(x_M^* - \frac{\Delta x}{2}, x_M^* + \frac{\Delta x}{2}\right)$$

2) Cuento el número de valores h_k de x que caen en cada intervalo centrado en x_k^* .

Estos dos pasos los hacen las funciones:

[n,bin] = hist(X,M): dado un número de intervalos M indica cuántos elementos del vector X caen en cada intervalo centrado en los valores del vector bin .

[n,bin] = histc(X,edges): indica cuántos elementos del vector X caen en cada intervalo limitado por los valores del vector edges .

Histograma de frecuencias (II)

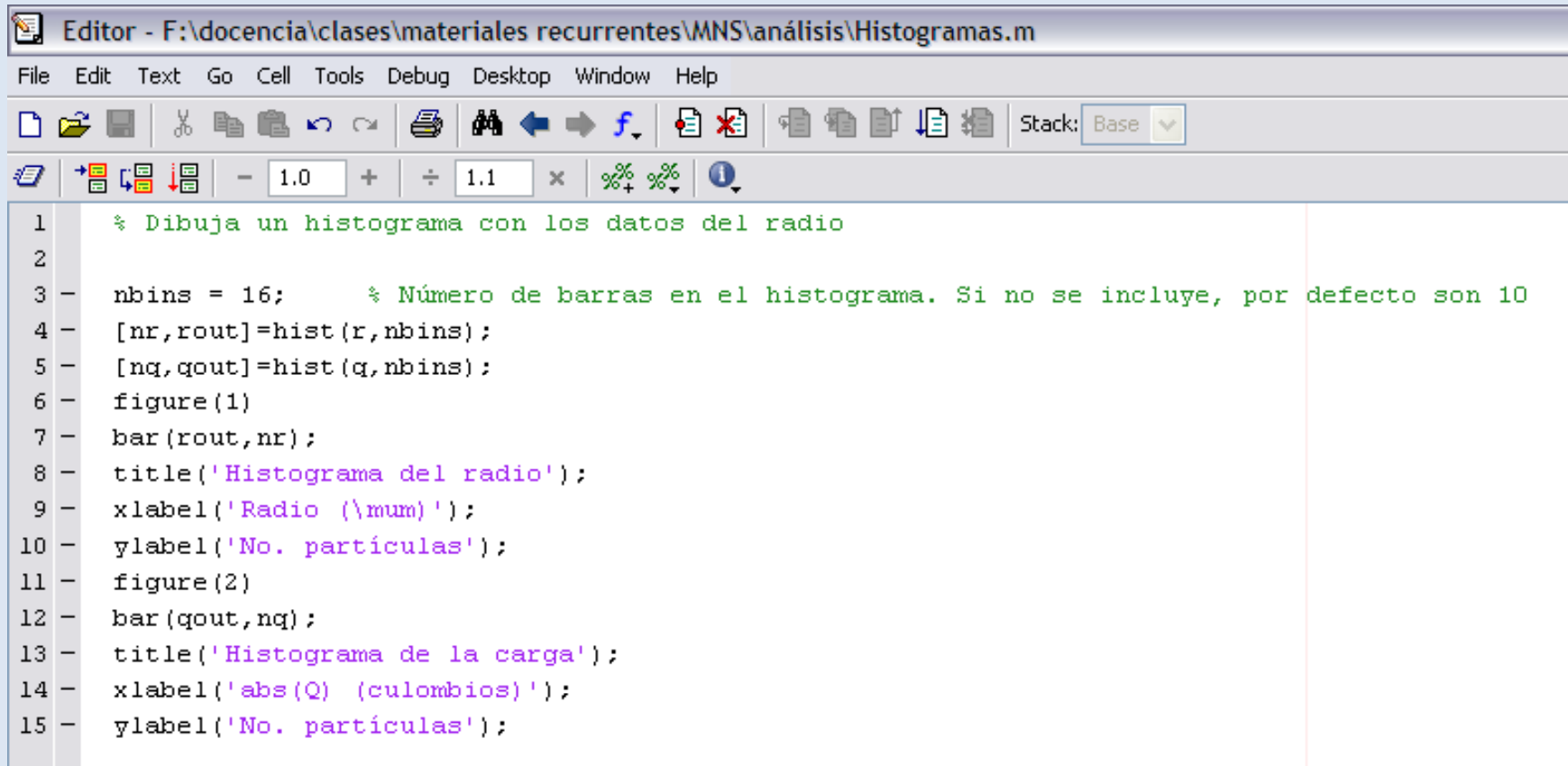
3) Represento h_k frente a x_k^*

De eso se encargan las funciones:

hist(X,bin): dibuja el histograma de los valores de X que caen en cada uno de los intervalos centrados en los valores del vector bin.

bar(n,bin): dibuja el histograma de los valores del vector n frente a los valores del vector bin.

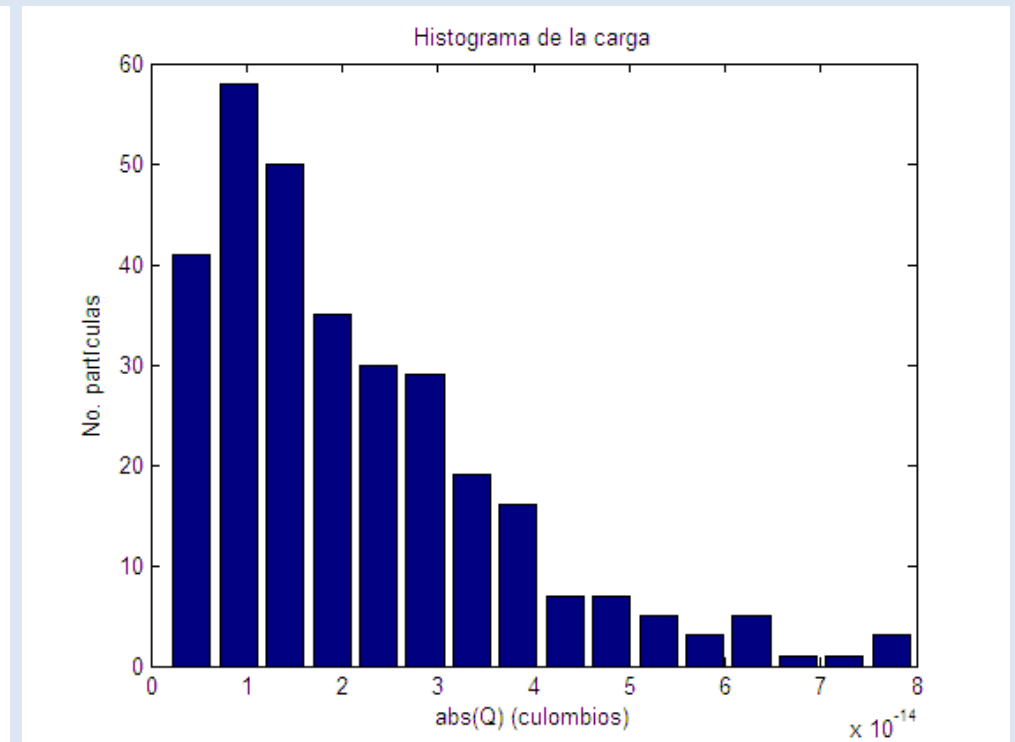
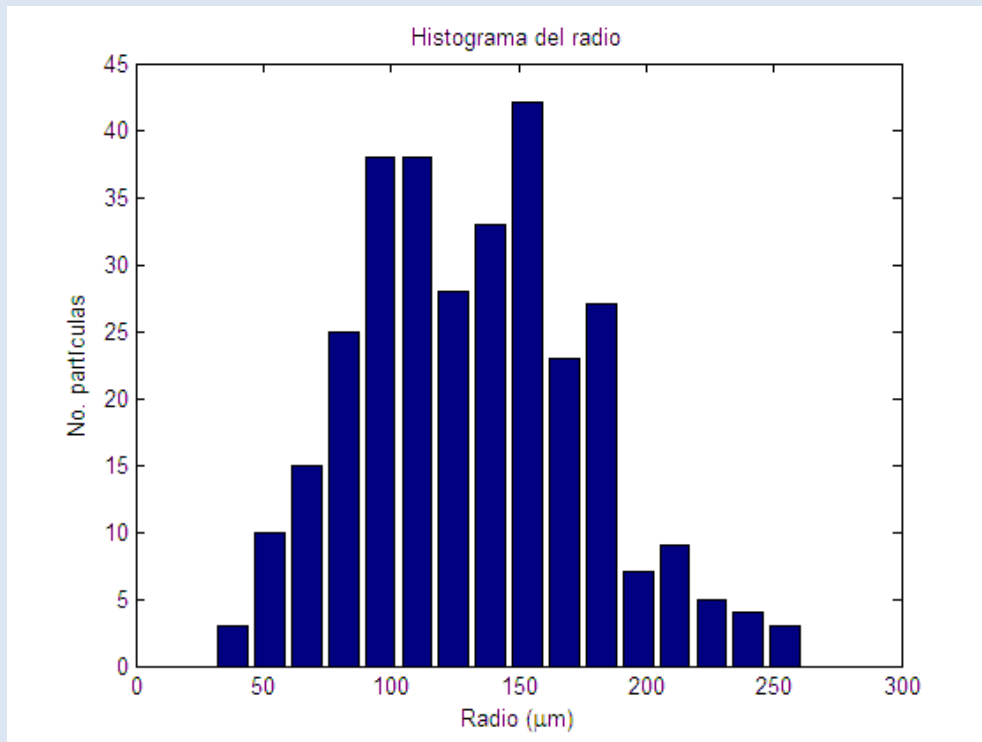
Ejemplo (histograma)



The image shows a screenshot of a MATLAB editor window. The title bar reads "Editor - F:\docencia\clases\materiales recurrentes\MNS\análisis\Histogramas.m". The menu bar includes "File", "Edit", "Text", "Go", "Cell", "Tools", "Debug", "Desktop", "Window", and "Help". The toolbar contains various icons for file operations, editing, and execution. Below the toolbar, there are zoom controls showing "1.0" and "1.1" magnification levels. The main editing area contains the following MATLAB code:

```
1 % Dibuja un histograma con los datos del radio
2
3 - nbins = 16; % Número de barras en el histograma. Si no se incluye, por defecto son 10
4 - [nr,rout]=hist(r,nbins);
5 - [nq,qout]=hist(q,nbins);
6 - figure(1)
7 - bar(rout,nr);
8 - title('Histograma del radio');
9 - xlabel('Radio (\mum)');
10 - ylabel('No. partículas');
11 - figure(2)
12 - bar(qout,nq);
13 - title('Histograma de la carga');
14 - xlabel('abs(Q) (culombios)');
15 - ylabel('No. partículas');
```

Ejemplo (histogramas)



Una pregunta recurrente es: ¿cuántas barras es razonable poner?

Volveremos a considerar esta pregunta más adelante

Distribuciones diferenciales a partir del histograma

El histograma no es lo mismo que la función de distribución diferencial.

Recuerda que:
$$f(x_o) = \left. \frac{dF}{dx} \right|_{x=x_o} \approx \frac{F(x_o + \Delta x/2) - F(x_o - \Delta x/2)}{\Delta x}$$

El número de datos x_j comprendidos entre los valores $x_o^* + \Delta x/2$ y $x_o^* - \Delta x/2$ viene dado por:

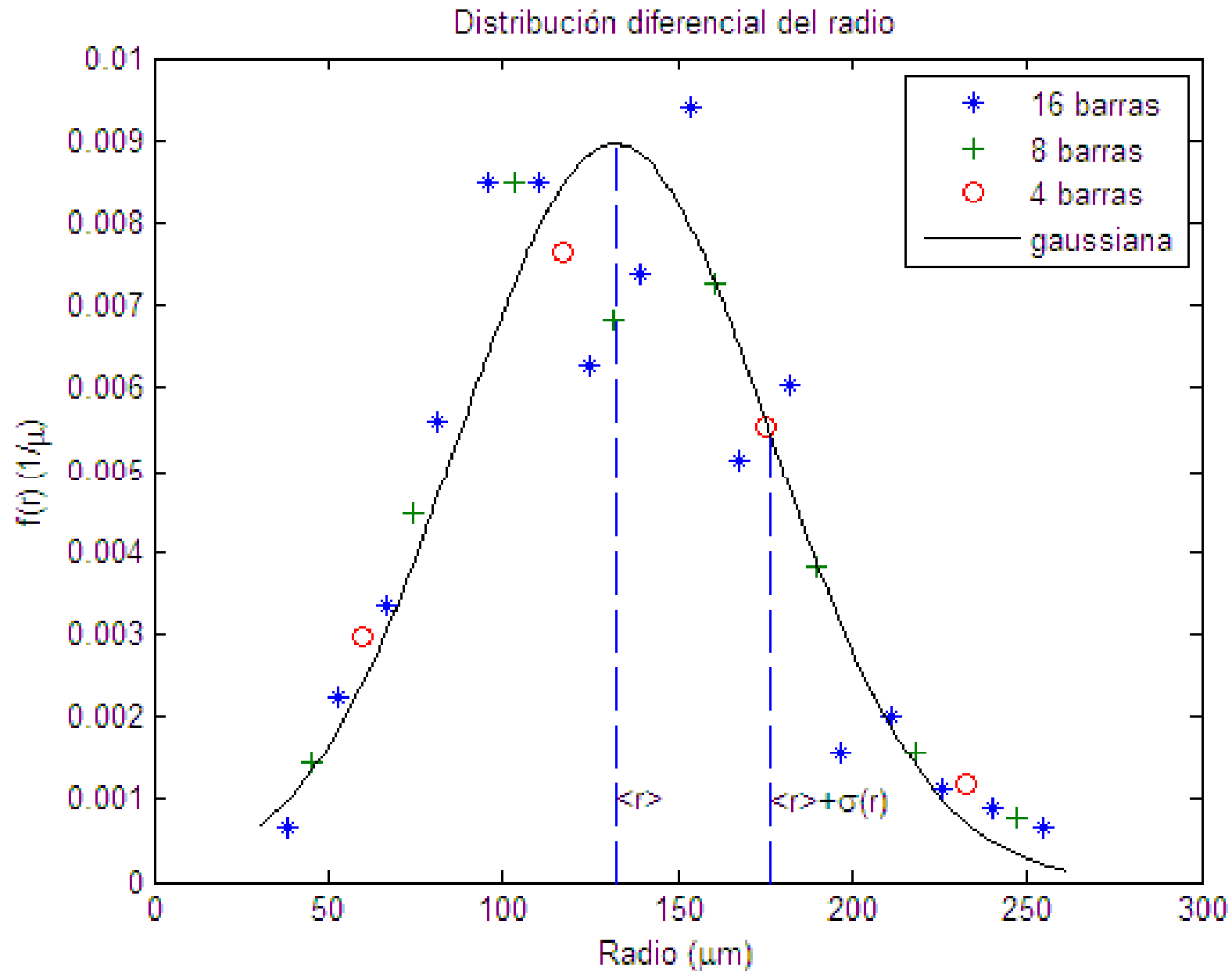
$$N \left[F(x_o^* + \Delta x/2) - F(x_o^* - \Delta x/2) \right]$$

Eso significa que si queremos estimar la distribución diferencial $f_{\text{exp}}(x)$, tenemos que dividir cada columna del histograma por su anchura y por el número total de datos.

Ejemplo (distribución diferencial)

```
Editor - F:\docencia\clases\materiales recurrentes\MNS\análisis\EjemploDiffR.m
File Edit Text Go Cell Tools Debug Desktop Window Help
[Icons] Stack: Base
- 1.0 + ÷ 1.1 x % % !
1  % Estima la distribución diferencial con los datos del radio
2  % a partir del histograma usando distinto número de intervalos.
3
4 - r = [r1 ; r2 ; r3 ; r4 ; r5];           % Concateno todas las series del radio
5 - sr = std(r,0);
6 - mr = mean(r);
7 - R = sort(r);
8 - [N M] = size(r);
9 - rmin= min(r);
10 - rmax= max(r);
11 - [n16,rout16]=hist(r,16);              % Datos para el histograma usando 16 barras
12 - [n8,rout8]=hist(r,8);                 % Datos para el histograma usando 8 barras
13 - [n4,rout4]=hist(r,4);                 % Datos para el histograma usando 4 barras
14 - dr16 = rout16(2)-rout16(1);
15 - fr16 = n16/(N*dr16);                  % Estimación de la densidad de probabilidad usando 16 barras
16 - dr8 = rout8(2)-rout8(1);
17 - fr8 = n8/(N*dr8);                    % Estimación de la densidad de probabilidad usando 16 barras
18 - dr4 = rout4(2)-rout4(1);
19 - fr4 = n4/(N*dr4);                    % Estimación de la densidad de probabilidad usando 16 barras
20 - Gdist = Gauss(R,mr,sr);
21 - figure(1);
22 - plot(rout16,fr16,'+',rout8,fr8,'+',rout4,fr4,'o',R,Gdist,'k-');
23 - title('Distribución diferencial del radio');
24 - xlabel('Radio (\mum)');
25 - ylabel('f(r)');
26 - legend('16 barras','8 barras','4 barras','gaussiana');
27
```

Ejemplo (distribución gaussiana)



La distribución gaussiana

- La distribución que hemos pintado en la figura anterior recibe el nombre de distribución gaussiana o distribución normal.

- Su ecuación es:

$$f_{gauss}(x; \langle x \rangle, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x - \langle x \rangle)^2}{2\sigma^2}\right]$$

σ : desviación estándar de la variable x

- Y para la distribución acumulativa:

$$F_{gauss}(x; \langle x \rangle, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x \exp\left[-\frac{(x_o - \langle x \rangle)^2}{2\sigma^2}\right] dx_o$$

Gausiana cumulativa: Implementación en matlab

Matlab tiene las siguientes funciones

$y = \mathbf{erf}(x)$: función de error.

$$y = \frac{2}{\sqrt{\pi}} \int_0^x \exp(-t^2) dt$$

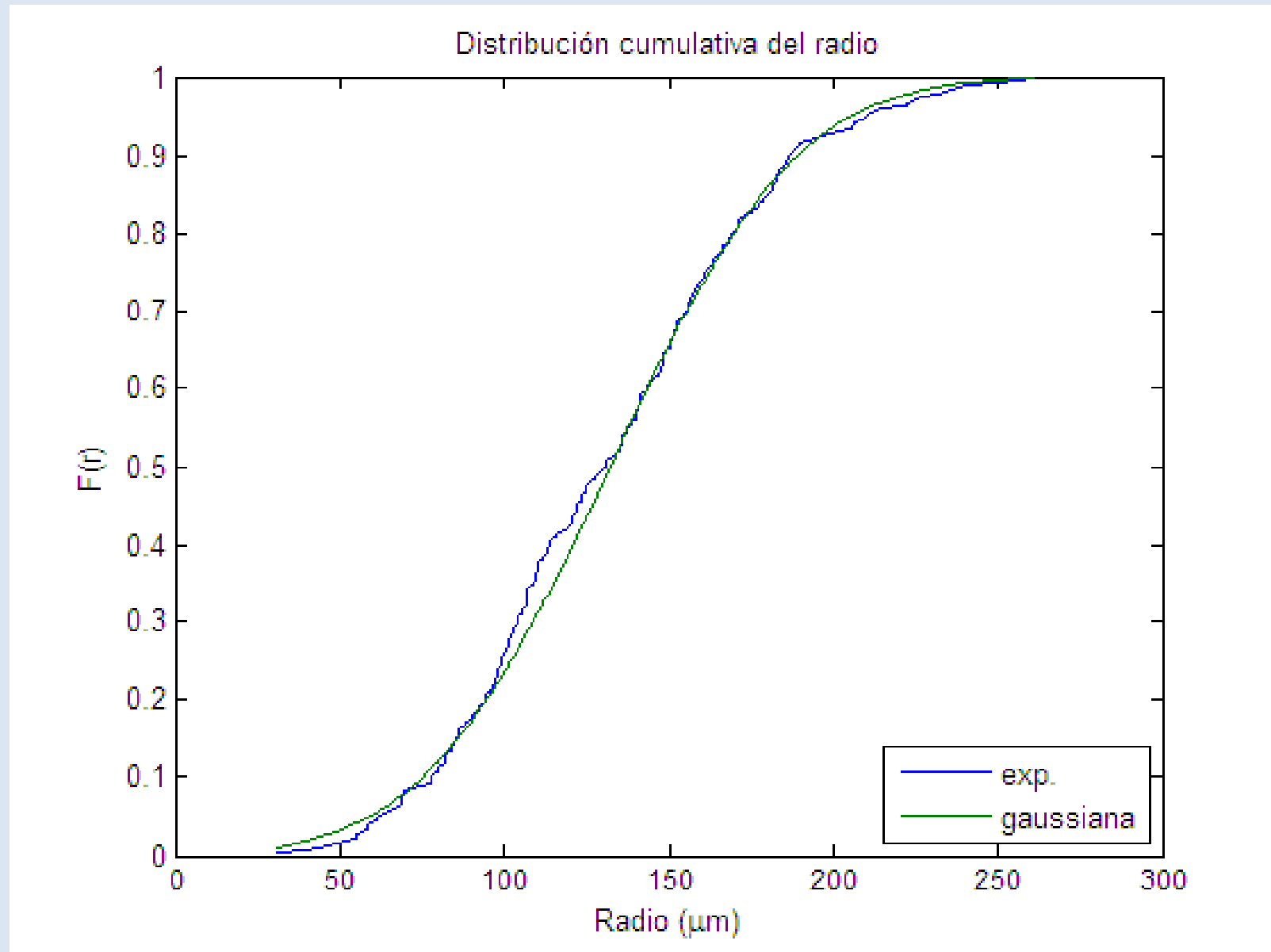
$y = \mathbf{erfc}(x)$: función de error complementaria

$$y = \frac{2}{\sqrt{\pi}} \int_x^{+\infty} \exp(-t^2) dt = 1 - \mathbf{erf}(x)$$

Para la distribución gausiana cumulativa:

$$F_{gauss}(x; \langle x \rangle, \sigma) = 1 - 0.5 \mathbf{erfc}\left(\frac{z}{\sqrt{2}}\right) \quad z = \frac{x - \langle x \rangle}{\sigma}$$

Ejemplo (gausiana acumulativa)



¿Cuántos intervalos escogemos?

- Cuando queremos representar un histograma, se nos presenta un dilema:
 - Con muchos intervalos, la distribución tiene muchos picos espúreos.
 - Con pocos intervalos, es difícil ver la forma de $f_{\text{exp}}(x)$.
- Si $f_{\text{exp}}(x)$ es una distribución gaussiana, el valor óptimo del intervalo viene dado por:

$$\Delta x \approx 3.49 \frac{s(x)}{N^{1/3}}$$

$s(x)$: desviación estándar de la muestra.

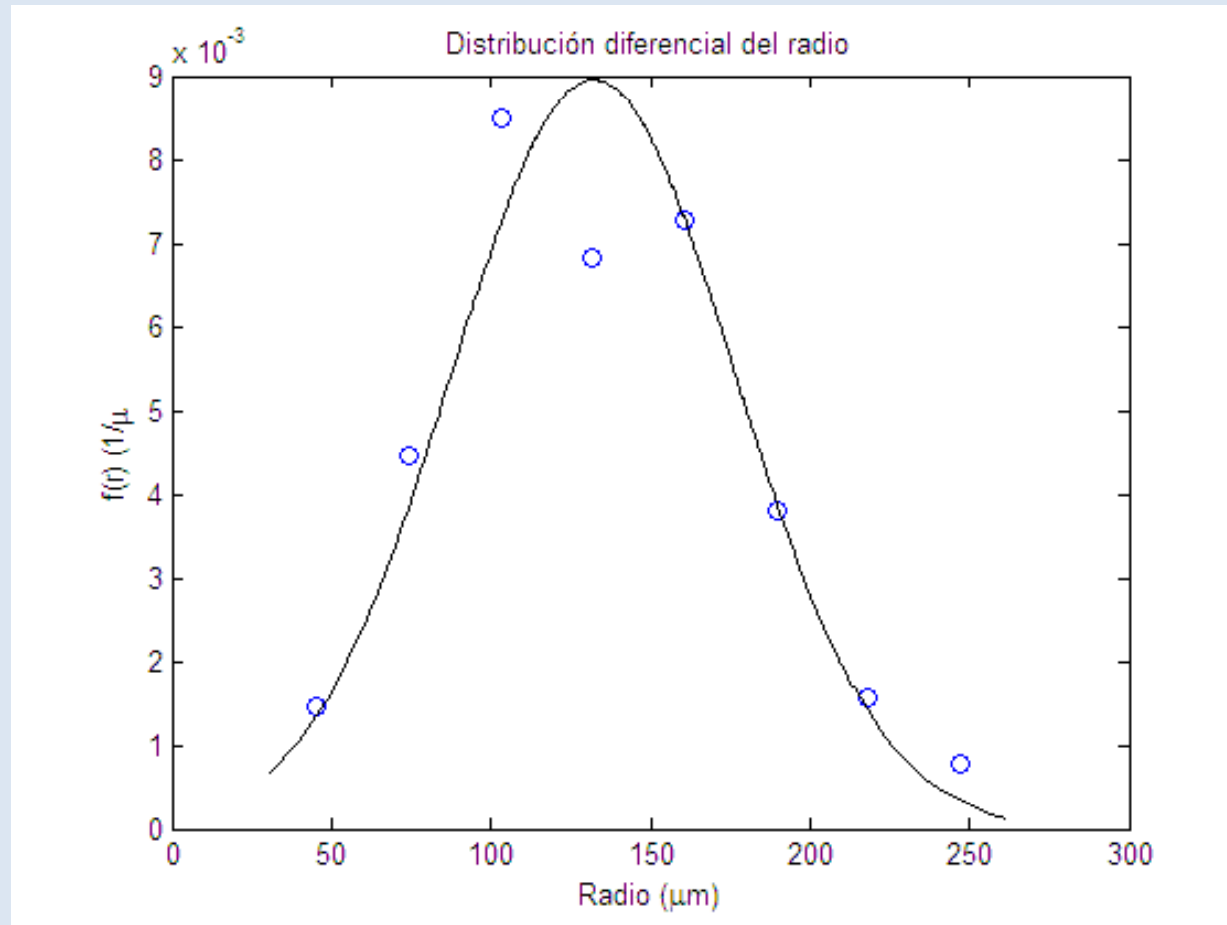
N : número de valores de la muestra.

- Si la distribución no es gaussiana, se puede usar esta regla como una primera estimación.

Ejemplo

En el caso de la distribución del radio, $N = 310$, $s(R) = 44.54 \mu\text{m}$

$$\Delta R \approx 3.49 \frac{s(R)}{N^{1/3}} = 3.49 \frac{44.54}{310^{1/3}} = 22.97 \mu\text{m}$$

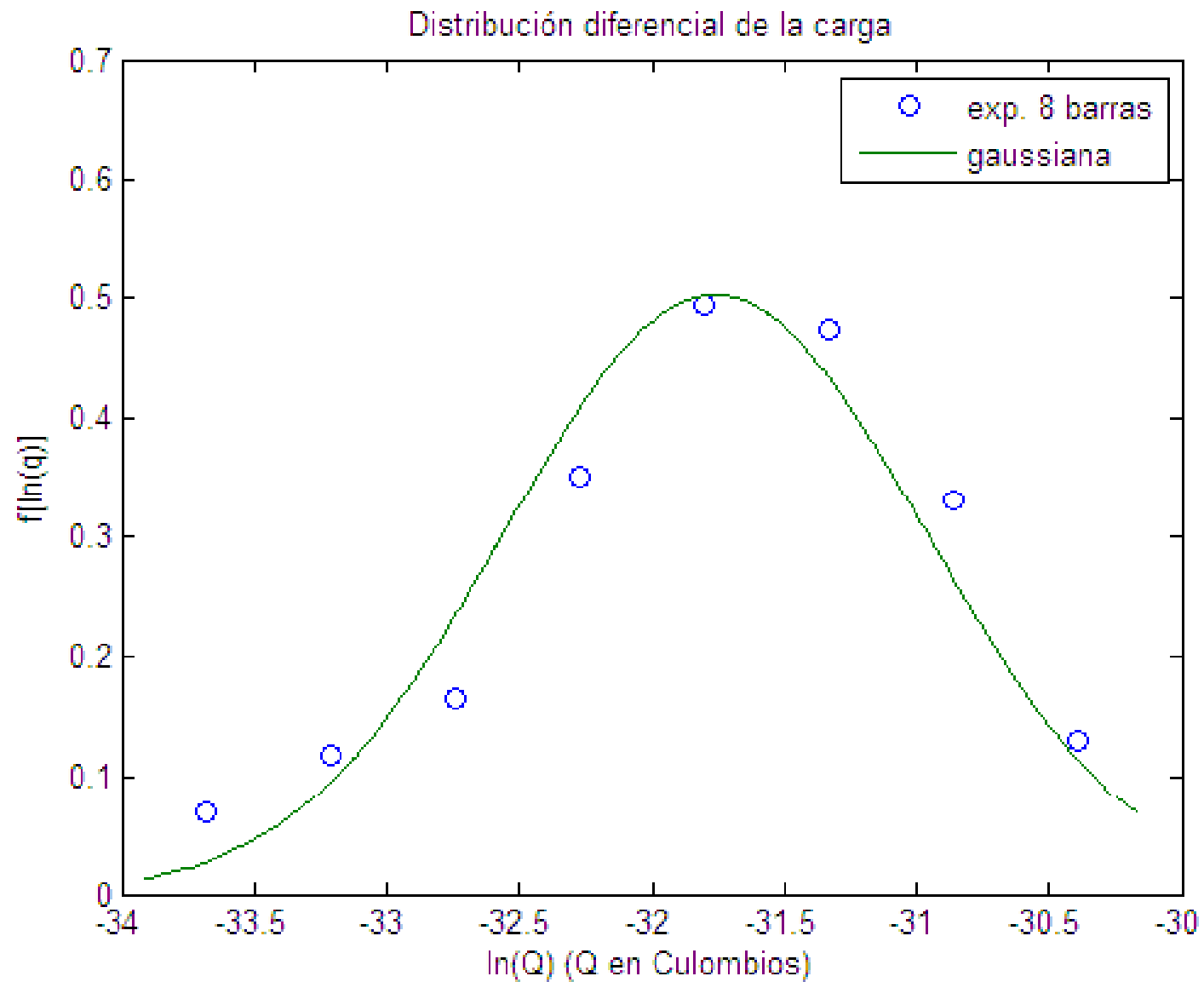


Ejemplo

Cuando calculamos la distribución diferencial para la carga, nos damos cuenta de que es el $\ln(Q)$ el que está distribuido (aproximadamente) como una distribución gaussiana.

```
Editor - F:\docencia\clases\materiales recurrentes\MNS\análisis\EjemploDiffLogQ.m
File Edit Text Go Cell Tools Debug Desktop Window Help
[Icons] Stack: Base
- 1.0 + 1.1 x % % !
1  % Dibuja la distribución diferencial del ln(q)
2
3  - q = [q1 ; q2 ; q3 ; q4 ; q5];      % Concateno todas las series de la carga
4  - Q = sort(q);
5  - [N M] = size(q);
6  - LogQ = log(Q);                    % Calculo el logaritmo neperiano de cada valor
7  - meanLQ = mean(LogQ);
8  - sLQ = std(LogQ,0);
9  - [n,bin]=hist(LogQ,8);             % Calculo el histograma de ln(q) con 8 barras
10 - d = bin(2)-bin(1);
11 - f = n/(N*d);                      % Estimación de la distribución diferencial
12 - Gdist = Gauss(LogQ,meanLQ,sLQ);
13 - plot(bin,f,'o',LogQ,Gdist);
14 - title('Distribución diferencial de la carga');
15 - xlabel('ln(Q) (Q en Culombios)');
16 - ylabel('f[ln(q)]');
17 - legend('exp. 8 barras','gaussiana')
```

Ejemplo



Cambios de variable en la distribución diferencial

Supongamos que queremos cambiar de la variable x a otra variable y que nos resulte más conveniente.

Para respetar la probabilidad, ha de cumplirse:

$$f(y)dy = f(x)dx \quad \rightarrow \quad f(y) = f(x) \frac{dx}{dy}$$

Por ejemplo, en el caso de la carga de las partículas queremos pasar de la carga q al logaritmo de la carga $\ln(q)$

Hacemos $y = \ln(q)$, $x = q$

$$\frac{dy}{dx} = \frac{d(\ln q)}{dq} = \frac{1}{q} \quad \rightarrow \quad \frac{dx}{dy} = q$$

$$f(\ln q) = f(q)q$$

La distribución log-normal

Sabemos que:

$$y = \ln(q) \quad f_{gauss}(y; \langle y \rangle, \sigma_y) = \frac{1}{\sigma_y \sqrt{2\pi}} \exp\left[-\frac{(y - \langle y \rangle)^2}{2\sigma_y^2}\right]$$

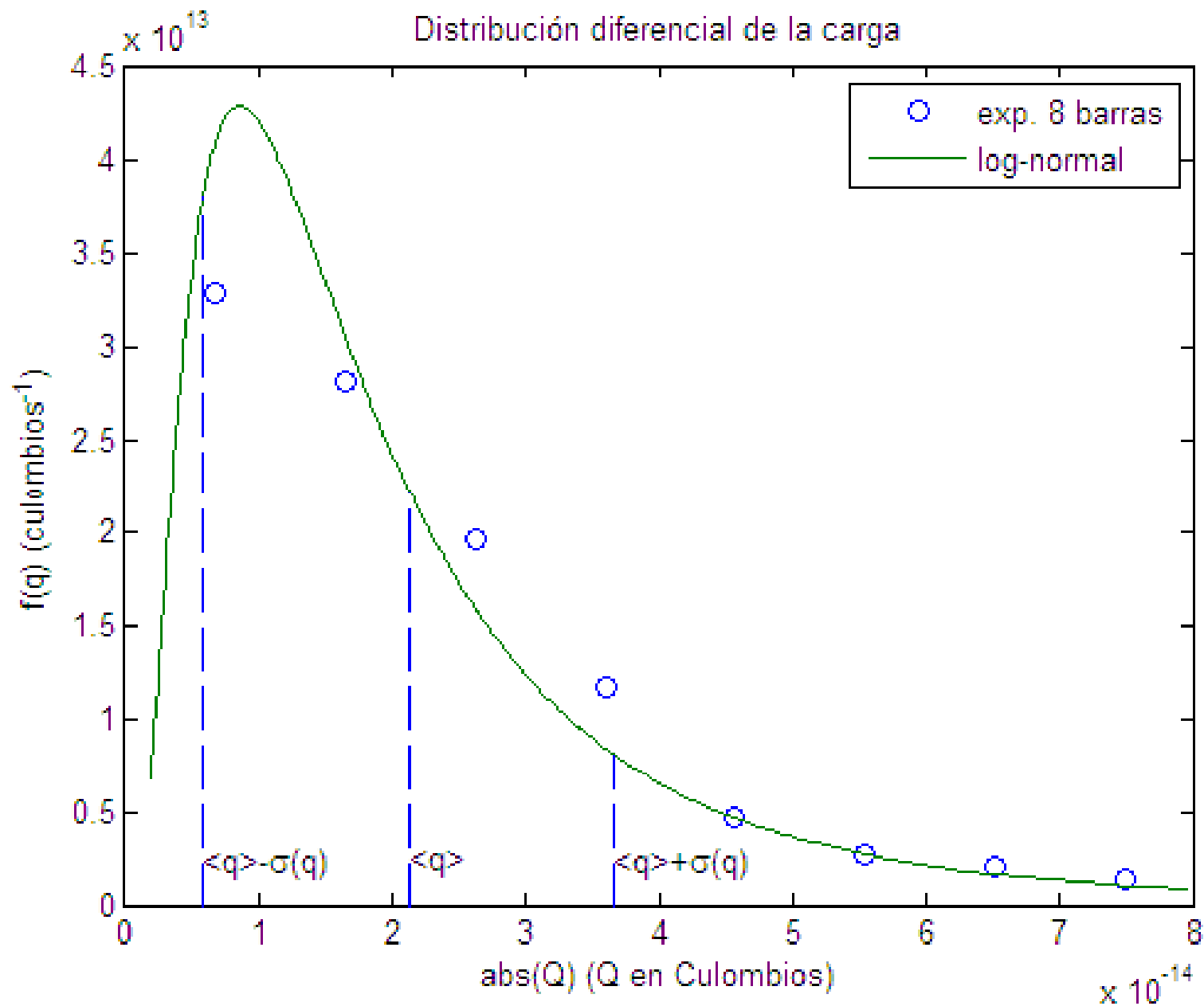
De lo que hemos dicho:

$$f(\ln q) = f(q)q \quad \Rightarrow \quad f(q) = \frac{1}{q} f(\ln(q))$$

$$f_{\log\text{-norm}}(q) = \frac{1}{q} \frac{1}{\sigma_{\ln(q)} \sqrt{2\pi}} \exp\left[-\frac{(\ln q - \langle \ln q \rangle)^2}{2\sigma_{\ln q}^2}\right]$$

A esta distribución se la llama distribución log-normal

Ejemplo (distribución log-normal)



Otras distribuciones

Otras distribuciones que trae Matlab y que nos resultarán útiles son:

- La distribución beta-incompleta.
- La distribución gamma-incompleta

Matlab trae sólo las distribuciones cumulativas.

Las distribuciones diferenciales y la inversas de estas funciones de las distribuciones cumulativas están programadas en un paquete extra llamado *Statistics toolbox*.

La función beta

- La función beta se define como:

En matlab: $y = \mathbf{beta}(a,b)$.

$$B(a,b) = \int_0^1 t^{a-1} (1-t)^{b-1} dt$$

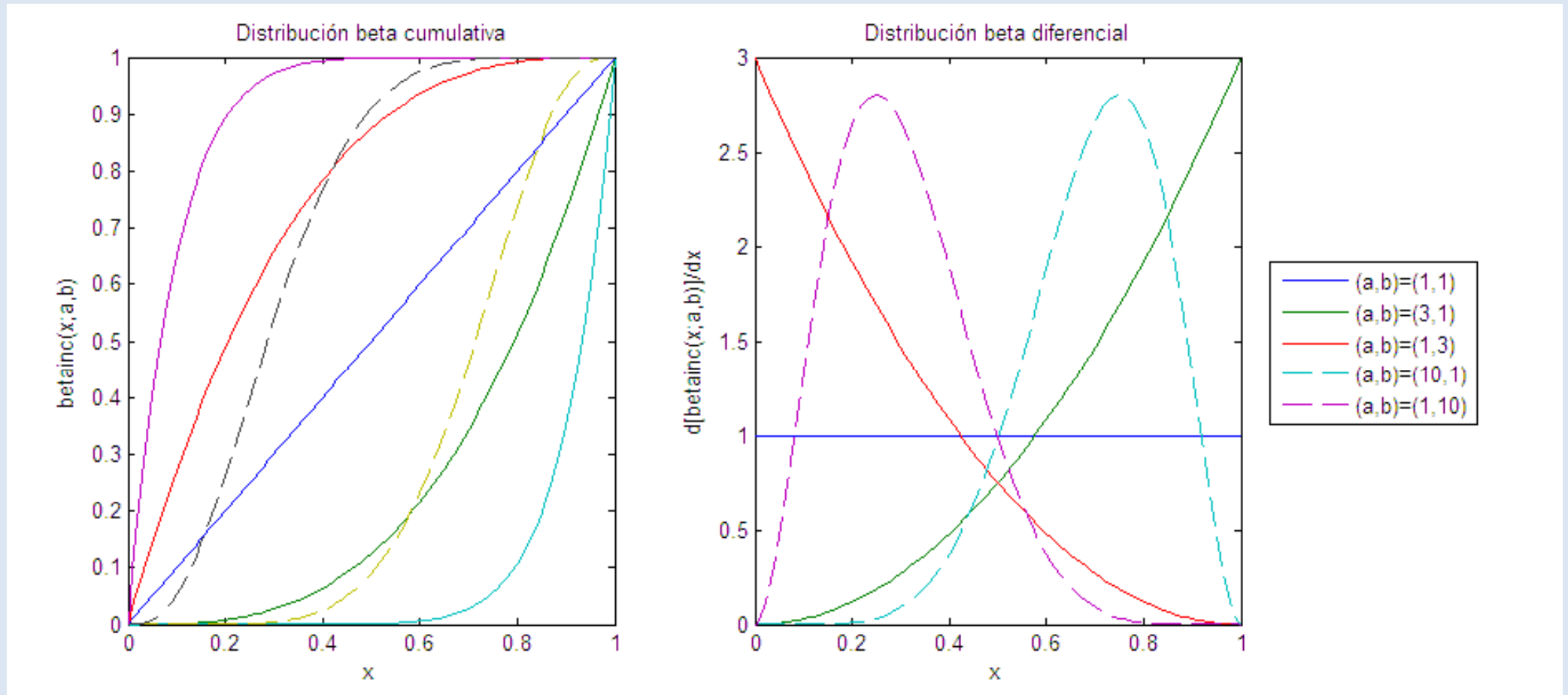
- La función beta incompleta es:

$$I_x(a,b) = \frac{1}{B(a,b)} \int_0^x t^{a-1} (1-t)^{b-1} dt$$

En matlab: $y = \mathbf{betainc}(x,a,b)$

Con $a, b > 0$ y $x < 1$.

Forma de la distribución beta



La función gamma

- Se define la función gamma como:

En matlab: $y = \mathbf{gamma}(a)$.

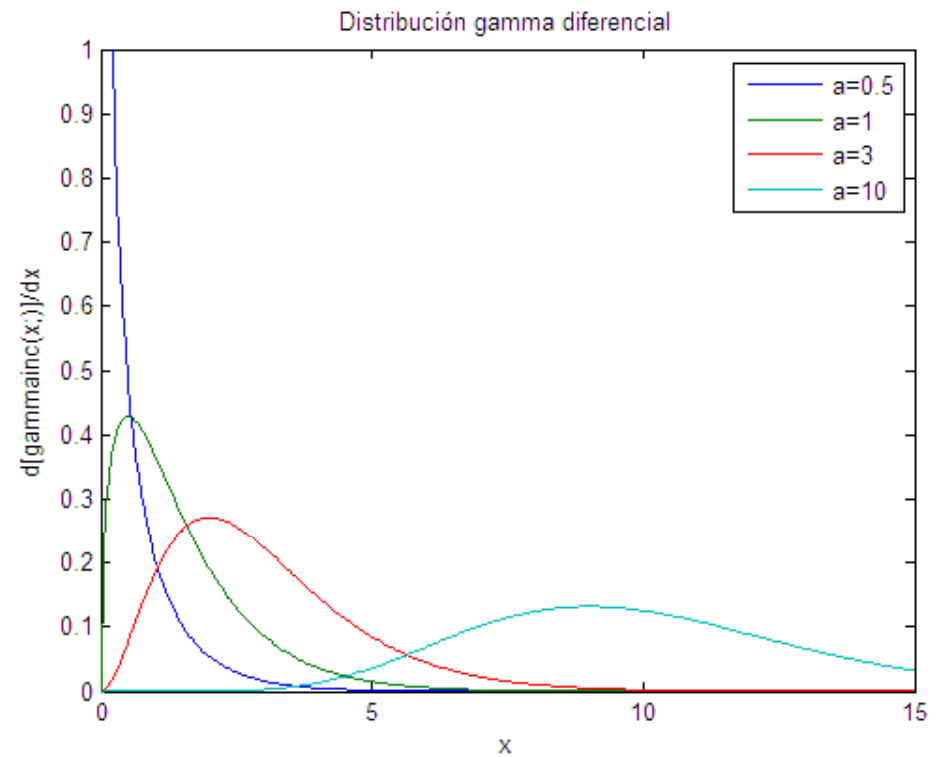
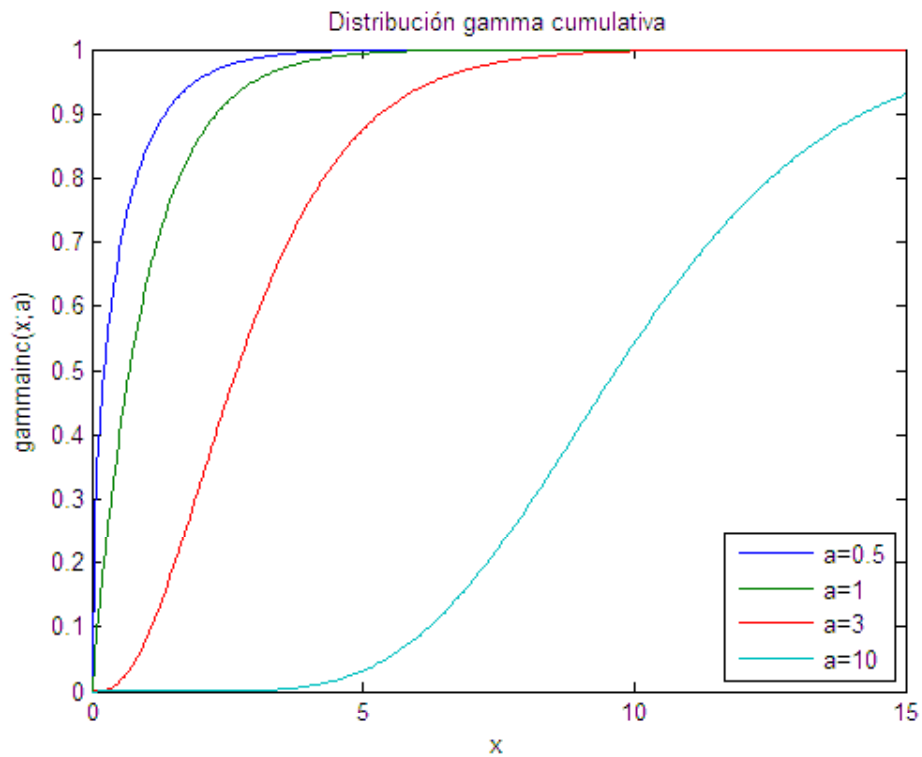
$$\Gamma(a) = \int_0^{\infty} e^{-t} t^{a-1} dt$$

- Se define la función gamma incompleta como:

En matlab: $y = \mathbf{gammainc}(x,a)$.

$$P(x, a) = \frac{1}{\Gamma(a)} \int_0^x e^{-t} t^{a-1} dt$$

Forma de la distribución gamma



Indeterminación en la media

- Supongamos que hemos tomado una muestra de N valores de x y que hemos encontrado que el valor medio de la muestra es $\langle x \rangle$.

Por ejemplo en nuestro caso, para el radio $\langle R \rangle = 132.16 \mu\text{m}$ para una muestra con $N=310$.

- El valor medio de x en la población, que llamaremos $\langle x \rangle_{\text{true}}$ nos es desconocido.

En nuestro caso, medir $\langle R \rangle_{\text{true}}$ supone medir el radio de todas las partículas del bote y hacer la media.

Se nos plantea la pregunta:

¿Está $\langle R \rangle$ suficientemente cerca de $\langle R \rangle_{\text{true}}$ o tenemos que medir más partículas?

Probabilidad de obtener un valor de la media

- Si cogemos M muestras, con el mismo número de valores N, la media $\langle x \rangle_k$ de cada una de ellas es una variable aleatoria. Como:

$$\langle x \rangle = \frac{1}{N} \sum_{i=1}^N x_i$$

- Si $P(x)$ es la probabilidad de obtener el valor x , entonces, la probabilidad de obtener la media $\langle x \rangle$ es:

$$P(\langle x \rangle) = P(x_1) \cdot P(x_2) \cdot \dots \cdot P(x_N)$$

condicionada a $x_1 + x_2 + x_3 + \dots + x_N = N \langle x \rangle$

Teorema central del límite

El Teorema Central del Límite nos dice que sea cual sea la forma de $P(x)$, cuando N es muy grande la distribución $P(\langle x \rangle)$ tiende a una distribución gaussiana con

$$\sigma(\langle x \rangle) = \frac{\sigma(x)}{\sqrt{N}}$$

$\sigma(x)$ es la desviación estándar de la población de x

Normalmente no conocemos la desviación estándar $\sigma(x)$ de la población, sino sólo la de una muestra, así que estimamos

$$\sigma(\langle x \rangle) \approx \frac{s(x)}{\sqrt{N}}$$

$s(x)$: desviación estándar de la muestra.

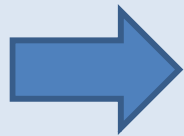
Podemos usar esta información para responder a la pregunta que nos hemos hecho.

Intervalo de confianza para la media

Como conocemos la forma de la distribución de $\langle x \rangle$, podemos dar la probabilidad de que la diferencia entre $\langle x \rangle$ y $\langle x \rangle_{true}$ sea mayor de un determinado valor.

Dado un número $k > 0$, la probabilidad P de que...

$$|\langle R \rangle_{true} - \langle R \rangle| < k\sigma(\langle R \rangle)$$



$$\langle R \rangle_{true} - k\sigma < \langle R \rangle < \langle R \rangle_{true} + k\sigma$$

...viene dada por:

$$P = \int_{\langle R \rangle_{true} - k\sigma}^{\langle R \rangle_{true} + k\sigma} f_{gauss}(t; \langle R \rangle_{true}, \sigma) dt = F_{gauss}(\langle R \rangle_{true} + k\sigma; \langle R \rangle_{true}, \sigma) - F_{gauss}(\langle R \rangle_{true} - k\sigma; \langle R \rangle_{true}, \sigma)$$

Cálculo del intervalo de confianza

Puede parecer que no podemos calcular P , porque no conocemos $\langle R \rangle_{true}$. Sin embargo, de la forma de F_{gauss} .

$$F_{gauss}(x; \langle x \rangle, \sigma) = 1 - 0.5 \operatorname{erfc}\left(\frac{z}{\sqrt{2}}\right) \quad z = \frac{x - \langle x \rangle}{\sigma}$$



$$F_{gauss}(\langle R \rangle_{true} + k\sigma; \langle R \rangle_{true}, \sigma) = 1 - 0.5 \operatorname{erfc}\left(\frac{k}{\sqrt{2}}\right)$$

$$F_{gauss}(\langle R \rangle_{true} - k\sigma; \langle R \rangle_{true}, \sigma) = 1 - 0.5 \operatorname{erfc}\left(-\frac{k}{\sqrt{2}}\right)$$

$$P = 0.5 \left[\operatorname{erfc}\left(-\frac{k}{\sqrt{2}}\right) + \operatorname{erfc}\left(\frac{k}{\sqrt{2}}\right) \right]$$

Ejemplo de cálculo (I)

Para la muestra del radio:

$$N = 310 \quad \langle R \rangle = 132.16 \mu\text{m} \quad s(R) = 44.54 \mu\text{m}$$

Sabemos que $\langle R \rangle$ tiene una distribución gaussiana centrada en $\langle R \rangle_{\text{true}}$ micras y con

$$\sigma(\langle R \rangle) \approx \frac{44.54}{\sqrt{310}} = 2.53 \mu\text{m}$$

¿Cuál es la probabilidad de que $|\langle R \rangle_{\text{true}} - \langle R \rangle| < 2 \mu\text{m}$?

Respuesta:

$$k = \frac{2}{2.53} = 0.79$$

$$P = 0.5 \left[\text{erfc} \left(-\frac{0.79}{\sqrt{2}} \right) + \text{erfc} \left(\frac{0.79}{\sqrt{2}} \right) \right] = 0.57$$

Ejemplo de cálculo (II)

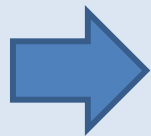
¿Cuántas partículas tengo que medir para que $|\langle R \rangle_{\text{true}} - \langle R \rangle| < 2 \mu\text{m}$ con probabilidad $P=0.95$?

Respuesta:

Necesito resolver la ecuación:

$$P = 0.5 \left[\text{erfc} \left(-\frac{k}{\sqrt{2}} \right) + \text{erfc} \left(\frac{k}{\sqrt{2}} \right) \right] = 0.95 \quad \rightarrow$$

$$k = 1.96 = \frac{2}{\sigma(\langle R \rangle)} \quad \text{Pero:} \quad \sigma(\langle R \rangle) \approx \frac{44.54}{\sqrt{N'}}$$



$$N' = \left(\frac{1.96}{2} \times 44.54 \right)^2 = 1905$$

Y estimo que tengo que agrandar mi muestra a $N' = 1905$ partículas.

¿Tienen dos series la misma media?

Supongamos que en un experimento he tomado una muestra 1 de una población y he hallado que la media es $\langle x \rangle_1$.

Si tomo otra muestra 2, en general, no obtendremos el mismo valor para la media, es decir:

$$\langle x \rangle_1 \neq \langle x \rangle_2$$

Se nos plantea la siguiente pregunta:

¿Es la diferencia fruto del azar, o ha cambiado la población entre las dos medidas?

La población puede haber cambiado porque:

- He alterado algún parámetro de mi experimento.
- Ha evolucionado en el tiempo transcurrido entre las dos medidas.

Ejemplo

Cada serie de datos de la carga y del radios se ha hecho con distintos parámetros experimentales. Por ejemplo, tenemos:

Serie 4

$$N = 72$$

$$\langle q \rangle = 2.04e-14 \text{ culombios}$$

$$s(q) = 1.24e-14 \text{ culombios}$$

$$\langle q \rangle = 0.146e-14 \text{ culombios}$$

Serie 5

$$N = 17$$

$$\langle q \rangle = 1.74e-14 \text{ culombios}$$

$$s(q) = 1.22e-14 \text{ culombios}$$

$$\sigma(\langle q \rangle) = 0.296e-14$$

La desviación estándar es parecida, pero la media es muy diferente. ¿Se debe la diferencia al azar? (p. ej. porque hemos tomado muy pocos datos en la serie 5). ¿Han afectado los cambios de los parámetros a los resultados del experimento?

El test de Student o t-test

- Este test determina si dos muestras A y B **con la misma desviación estándar** difieren de forma significativa en su media.

1) Calcular la desviación estándar conjunta s_D como:

$$s_D = \sqrt{\frac{(N_A - 1)s_A(x)^2 + (N_B - 1)s_B(x)^2}{N_A + N_B - 2} \left(\frac{1}{N_A} + \frac{1}{N_B} \right)}$$

2) Calcular el valor de:

$$t = \frac{\langle x \rangle_A - \langle x \rangle_B}{s_D}$$

De ahí el nombre de t-test.

- Llamemos $P(t|\nu)$ representa la probabilidad de obtener por azar un valor de $|t|$ mayor o igual que $|\langle x \rangle_A - \langle x \rangle_B|/s_D$

$\nu = N_A + N_B - 2$, recibe el nombre de número de grados de libertad.

El t-test (II)

Puede demostrarse que la probabilidad $P(t|\nu)$ viene dada por:

$$P(t|\nu) = \frac{1}{\nu^2 B\left(\frac{1}{2}, \frac{\nu}{2}\right)} \int_{-t}^t \left(1 + \frac{z^2}{\nu}\right)^{-\frac{\nu+1}{2}} dz$$

$P(t|\nu)$ puede escribirse usando la función beta incompleta.

$$P(t|\nu) = \text{betainc}\left(x, \frac{\nu}{2}, \frac{1}{2}\right) \quad x = \frac{\nu}{\nu + t^2}$$

A esta distribución también se la llama distribución de Student acumulativa.

Si $P(t|\nu)$ es bajo, entonces la diferencia entre las medias de las muestras A y B es poco probable que se deba al azar.

Se concluye entonces que las medias son diferentes.

Ejemplo (test-t)

Las cuentas de nuestro ejemplo son:

$$s_D = \sqrt{\frac{71 \times (1.24e-14)^2 + 16 \times (1.22e-14)^2}{87} \left(\frac{1}{72} + \frac{1}{17} \right)}$$

$$s_D = 3.33e-15$$

$$t = \frac{(2.04 - 1.74)e-14}{3.33e-15} = 0.901$$

$$x = \frac{87}{87 + 0.901^2} = 0.991$$

$$P(0.901 | 87) = \text{betainc} \left(0.991, \frac{87}{2}, \frac{1}{2} \right) = 0.376$$

La probabilidad de que la diferencia se deba al azar es de un 37.6%. Se concluye que es muy probable que ambas medias sean diferentes por haber cambiado los parámetros experimentales

Indeterminación en la desv. estándar

- Supongamos que hemos tomado una muestra de N valores de x y que la desviación estándar de la muestra es $s(x)$.

Por ejemplo en nuestro caso, para el radio $s(R) = 44.54 \mu\text{m}$ para una muestra con $N = 310$.

- La desviación estándar de la población, que llamaremos $\sigma(x)$ nos es desconocida.

En nuestro caso, medir $\sigma(x)$ supone medir el radio de todas las partículas del bote y hallar su desviación estándar.

Se nos plantea la pregunta:

¿Está $s(R)$ suficientemente cerca de $\sigma(R)$ o tenemos que medir más partículas?

NOTA: Este problema no es tan común como el de la media porque muchas veces no conocemos la forma de la distribución y/o no nos interesa saber con mucha precisión su desviación estándar.

Probabilidad de obtener una variancia

- Si tomamos M muestras, la variancia $\text{Var}_k(x)$ de cada una de ellas es una variable aleatoria. Como:

$$\text{Var}(x) = \frac{1}{N} \sum_{i=1}^N (x_i - \langle x \rangle)^2 = \langle x^2 \rangle - \langle x \rangle^2$$

- Si $P(x)$ es la probabilidad de obtener el valor x , entonces, la probabilidad de obtener $\text{Var}(x)$ es:

$$P(\text{Var}(x)) = P(x_1) \cdot P(x_2) \cdot \dots \cdot P(x_N)$$

condicionada a:

$$x_1^2 + x_2^2 + \dots + x_N^2 - N \langle x \rangle^2 = N \text{Var}(x)$$

$$x_1 + x_2 + \dots + x_N = \langle x \rangle$$

La distribución chi-cuadrado

Puede demostrarse teóricamente que si x_j es una variable aleatoria **con distribución gaussiana** con desviación estándar σ , la variable

$$\chi^2 = \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \langle x \rangle)^2 = \frac{(N-1)s(x)}{\sigma(x)}$$

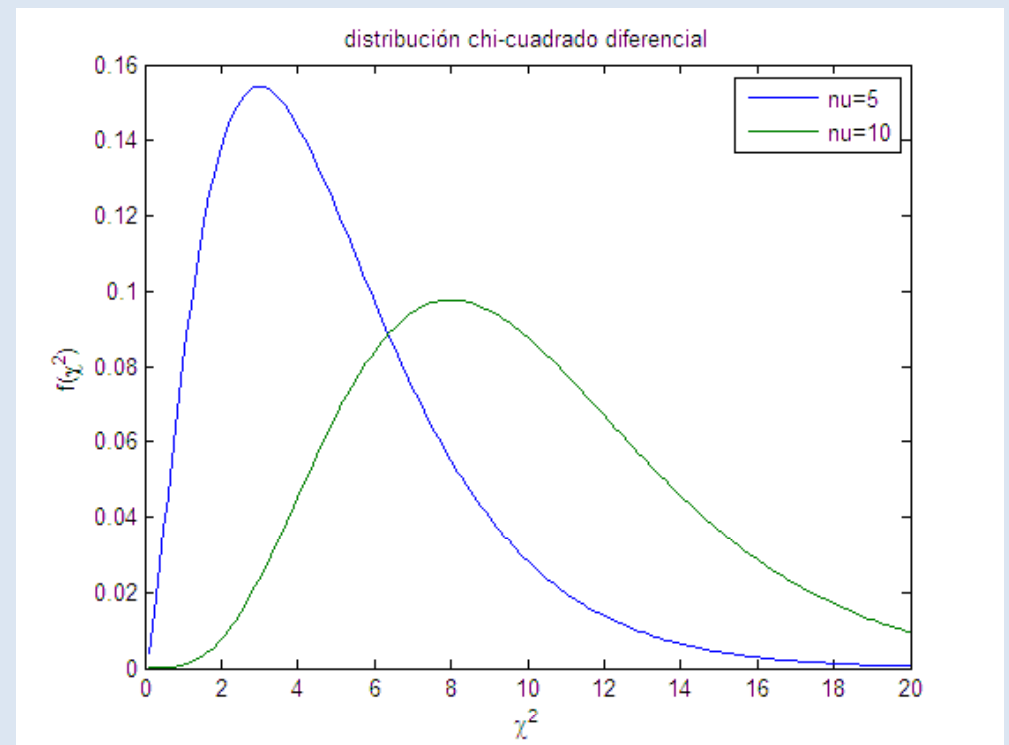
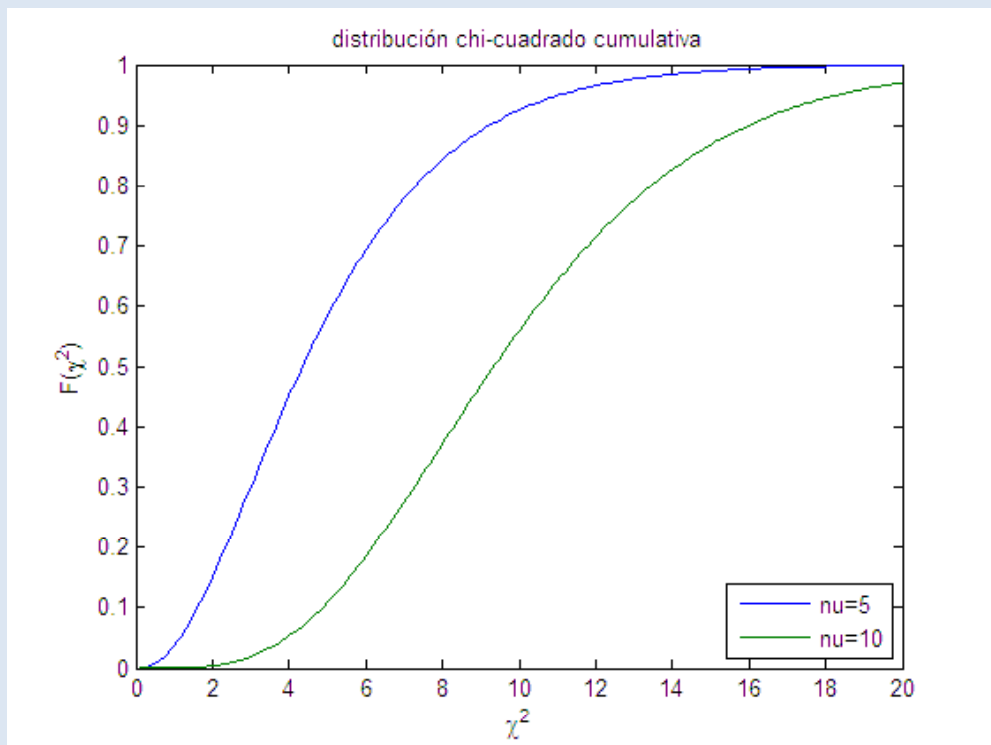
sigue la siguiente distribución acumulativa:

$$F_{\text{chi-cuad}}(\chi^2; \nu) = \frac{1}{2^{\frac{\nu}{2}} \Gamma\left(\frac{\nu}{2}\right)} \int_0^{\chi^2} t^{\frac{\nu}{2}-1} e^{-\frac{t}{2}} dt = \text{gammainc}\left(\frac{\chi^2}{2}; \frac{\nu}{2}\right)$$

que recibe el nombre de distribución chi-cuadrado

$\nu = N-1$ es el número de grados de libertad.

Representación gráfica (chi-cuadrado)



$$\langle \chi^2 \rangle = \nu, \text{Var}(\chi^2) = 2\nu$$

Intervalos de confianza para la varianza

Si **estamos seguros de que x sigue una distribución gaussiana**, podemos dar la probabilidad de que la diferencia entre $s(x)$ y $\sigma(x)$ sea mayor que un determinado valor.

Dado un número $k > 0$ queremos saber cuál es la probabilidad P de que:

$$|\sigma(R) - s(R)| < k\sigma(R) \quad \rightarrow$$

$$\sigma(R) - k\sigma(R) < s(R) < \sigma(R) + k\sigma(R)$$

$$(1 - k)\sigma(R) < s(R) < (1 + k)\sigma(R)$$

Multiplicando por $(N-1)$ y dividiendo por $\sigma(x)$

$$(1 - k)(N - 1) < \frac{(N - 1)s(R)}{\sigma(R)} < (1 + k)(N - 1)$$

Cálculo del intervalo de confianza

Y usando la definición de χ^2

$$(1-k)(N-1) < \chi^2 < (1+k)(N-1)$$

Luego la probabilidad P se encuentra partir de la distribución chi-cuadrado con $\nu=N-1$ grados de libertad haciendo:

$$P = \int_{(1-k)(N-1)}^{(1+k)(N-1)} f_{chi-cuad}(t; \nu) dt$$

$$P = F_{chi-cuad}((1+k)(N-1); \nu) - F_{chi-cuad}((1-k)(N-1); \nu)$$

Y usando la definición de la distribución chi-cuadrado:

$$P = \text{gammainc}\left(\frac{(1+k)\nu}{2}; \frac{\nu}{2}\right) - \text{gammainc}\left(\frac{(1-k)\nu}{2}; \frac{\nu}{2}\right)$$

$\nu=N-1$ es el número de grados de libertad.

Cálculo del intervalo de confianza

Sabemos que la distribución de R en la muestra de los datos del radio **se aproxima bien por una forma gaussiana**.

$$N = 310 \quad \langle R \rangle = 132.16 \mu\text{m} \quad s(R) = 44.54 \mu\text{m}$$

$$\text{Var}(R) = 1984.44 \text{ micras}^2$$

¿Cuál es la probabilidad de que $|\sigma - s(R)|$ sea menor que un 1% de $\sigma(R)$?

Respuesta: $k = 0.01$

$$P = \text{gammainc}\left(\frac{1.01 \times 309}{2}; \frac{309}{2}\right) - \text{gammainc}\left(\frac{0.99 \times 309}{2}; \frac{309}{2}\right)$$

$$P = 0.0976$$

¿Tienen dos series la misma variancia?

Supongamos que hemos tomado una muestra 1 de una población y hemos hallado que su variancia es $\text{Var}_1(x) = s_1(x)^2$.

Si sacamos otra muestra 2, en general, no obtendremos el mismo valor para la variancia, es decir:

$$s_1(x) \neq s_2(x)$$

Se nos plantea la siguiente pregunta:

¿Es la diferencia fruto del azar, o es que ha cambiado la población entre las dos medidas?

La población puede haber cambiado porque:

- He alterado algún parámetro de mi experimento.
- Ha evolucionado en el tiempo transcurrido entre las dos medidas.

¿Tienen dos series la misma variancia?

Ejemplo: si consideramos separadamente los valores del radio para la serie 4 y 5:

Serie 4:

$$N = 72$$

$$\langle R \rangle = 133.41 \mu\text{m}$$

$$s(R) = 40.64 \mu\text{m}$$

$$\sigma(\langle R \rangle) = 14.79 \mu\text{m}$$

Serie 5

$$N = 17$$

$$\langle R \rangle = 133.88 \mu\text{m}$$

$$s(R) = 34.29 \mu\text{m}$$

$$\sigma(\langle R \rangle) = 8.32 \mu\text{m}$$

La media es aproximadamente la misma, pero la desviación estándar es diferente. ¿Puede deberse la diferencia a que hemos tomado muy pocos datos en la serie 5?

Nota: Ambas series no tienen por qué tener la misma media para poder hacer el F-test

El test-F

1) Calcula el cociente

$$\text{Si } \text{Var}_A(x) > \text{Var}_B(x) \quad F = \text{Var}_A(x) / \text{Var}_B(x) \quad \nu_1 = N_A - 1 \quad \nu_2 = N_B - 1$$

$$\text{Si } \text{Var}_B(x) > \text{Var}_A(x) \quad F = \text{Var}_B(x) / \text{Var}_A(x) \quad \nu_1 = N_B - 1 \quad \nu_2 = N_A - 1$$

de modo que el valor de F sea siempre mayor que uno.

2) Llamo $P(F | \nu_1; \nu_2)$ a la probabilidad de obtener al azar un valor de F igual o mayor que el que he obtenido. Puede demostrarse que esa probabilidad vale:

$$P(F | \nu_1, \nu_2) = \int_0^F \frac{\Gamma\left(\frac{\nu_1 + \nu_2}{2}\right)}{\Gamma\left(\frac{\nu_1}{2}\right)\Gamma\left(\frac{\nu_2}{2}\right)} \left(\frac{\nu_1}{\nu_2}\right)^{\nu_1/2} \frac{z^{(\nu_1-2)/2}}{\left(1 + \frac{\nu_1}{\nu_2} z\right)^{\frac{\nu_1 + \nu_2}{2}}} dz$$

El test-F

3) Puede demostrarse que

$$I = \text{betainc}\left(\frac{v_2}{v_2 + v_1 F}; \frac{v_2}{2}, \frac{v_1}{2}\right)$$

$$P(F | v_1, v_2) = 2 \min[I, 1 - I]$$

Si $P(t|v)$ es bajo, entonces la diferencia entre las medias de las muestras A y B es poco probable que se deba al azar.

Se concluye entonces que las medias son diferentes.

Ejemplo (test-F)

Para la serie 4: $\text{Var}(R) = 1651.5$

Para la serie 5: $\text{Var}(R) = 1176.1$

$$F = \frac{1651.5}{1176.1} = 1.404 \quad N_1 = 72 \quad N_2 = 17$$

$$\frac{\nu_2}{\nu_2 + \nu_1 F} = \frac{16}{16 + 71 \times 1.404} = 0.1383$$

$$I = \text{betainc}\left(0.1383; \frac{16}{2}, \frac{71}{2}\right) = 0.2275$$

$$P(F | \nu_1, \nu_2) = 2 \min[I, 1 - I] = 0.455$$

Se concluye entonces que hay un 45.5% de posibilidades de que la desviación estándar de ambas muestras sea igual.

¿Se ajusta una distribución a una fórmula teórica?

Supongamos que hemos tomado una muestra de N valores de x y que nos han dicho que la expresión de la distribución de x en la población es $F(x)$.

Representamos $F_{\text{exp}}(x)$ y no coincide con $F(x)$

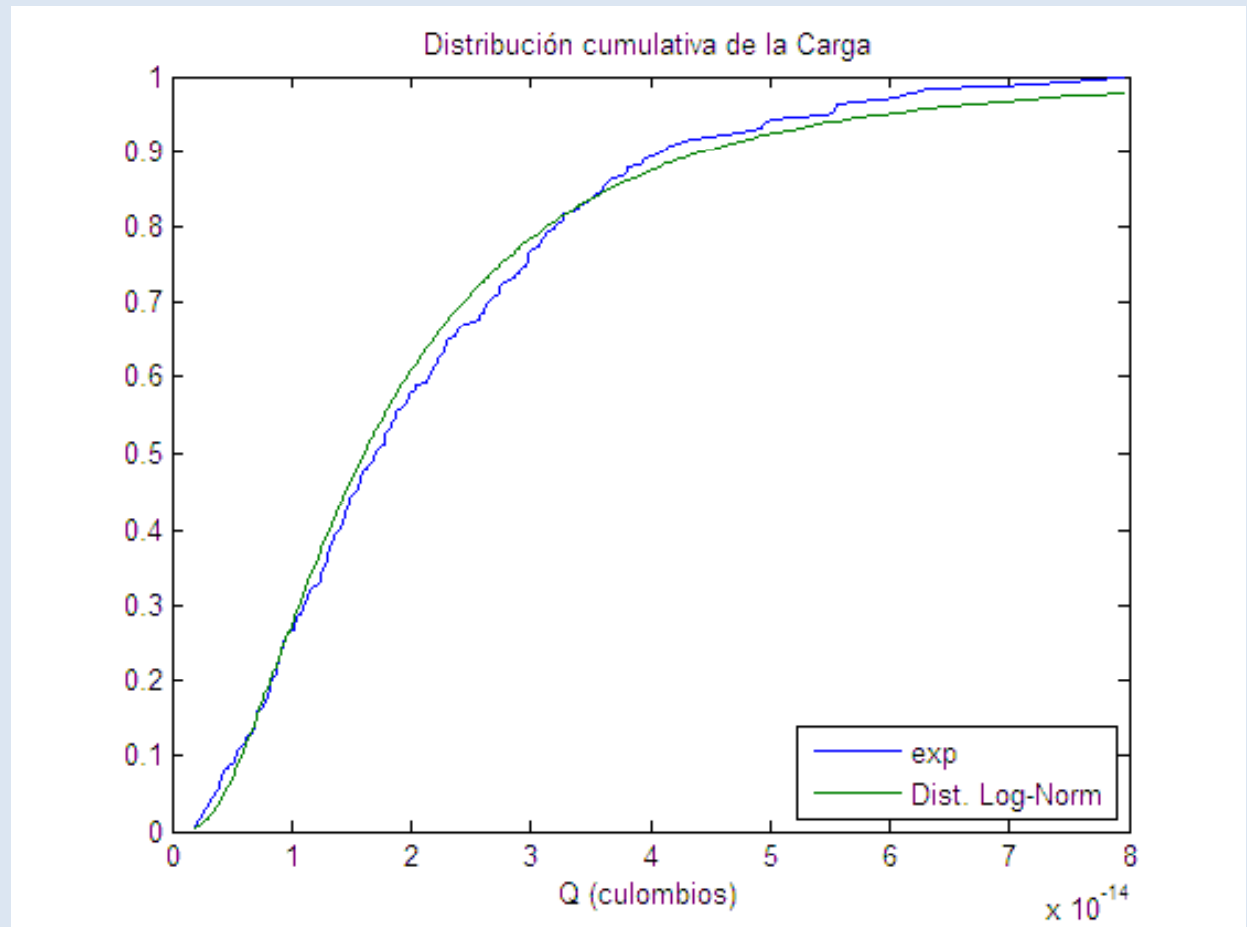
Se nos plantea la pregunta

¿Es esto fruto del azar porque hemos tomado muy pocos valores de N ? O por el contrario ¿es falso que la distribución de x en la población es $F(x)$?

Ejemplo

Por ejemplo, parece que los datos de la carga se ajustan a una distribución log-normal:

¿Cuál es la probabilidad de que las diferencias se deban al azar?



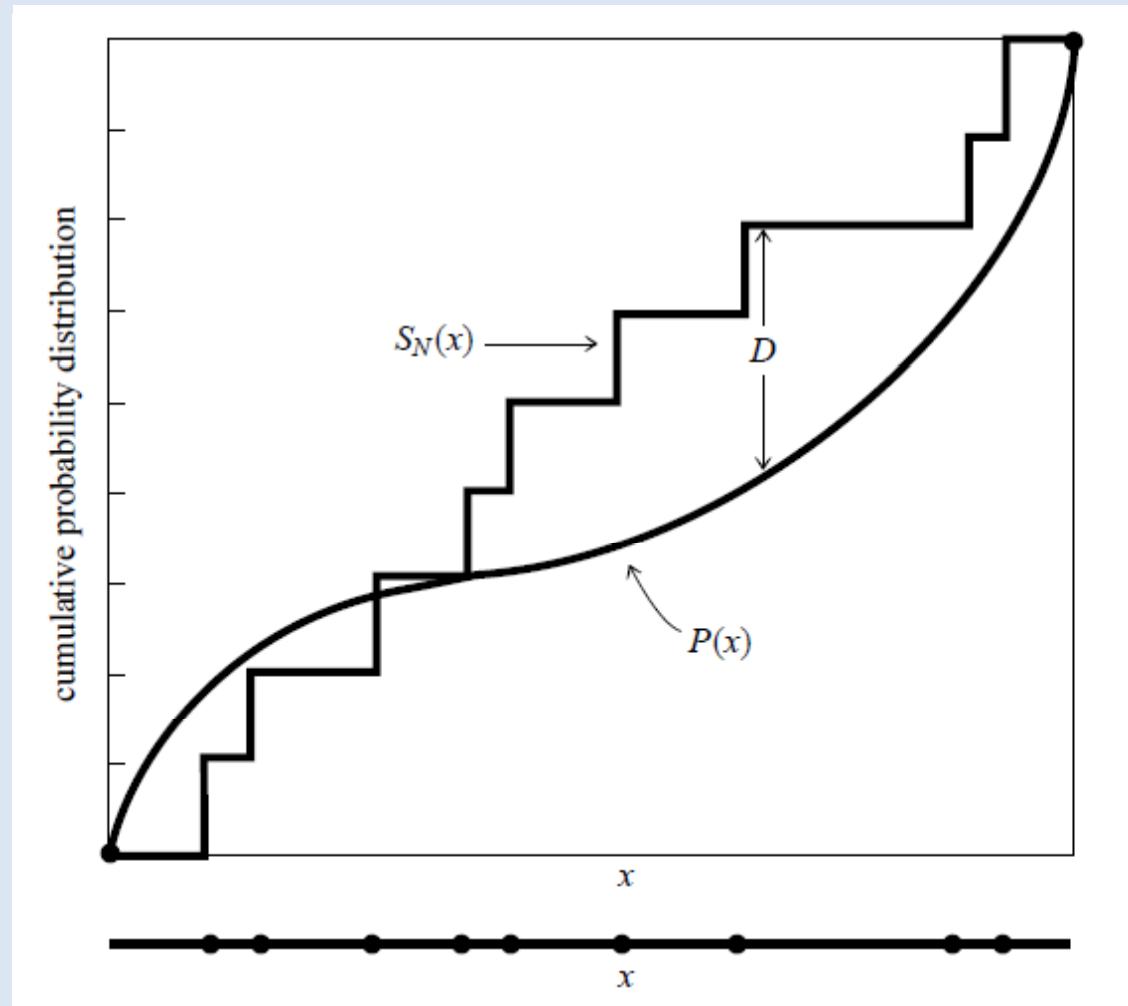
El test de Kolmogorov-Smirnoff

1) Calculamos:

$$D_o = \max |F_{\text{exp}}(x) - F_{\text{teo}}(x)|$$

$F_{\text{teo}}(x)$:
distribución
cumulativa teórica
que queremos
comprobar

D : máxima
distancia en
vertical entre la
distribución
cumulativa
experimental y la
teórica que
proponemos.



El test de Kolmogorov-Smirnoff

Puede demostrarse que la probabilidad de haber obtenido un valor de la distancia D mayor o igual que D_o por azar viene dada por:

$$P(D > D_o) = 2 \sum_{j=1}^{\infty} (-1)^{j-1} \exp(-2j^2 z^2)$$

$$z \approx \left[\sqrt{N} + 0.12 + \frac{0.11}{\sqrt{N}} \right] D_o$$

Si $P(D > D_o)$ es bajo, es muy probable que $F_{\text{exp}}(x)$ no se corresponda con $F_{\text{teo}}(x)$.

Código

```
Editor - F:\docencia\clases\materiales recurrentes\MNS\ analisis\KSone.m
File Edit Text Go Cell Tools Debug Desktop Window Help
[Icons] Stack: Base
Save - 1.0 + ÷ 1.1 x % %
1 function [prob,d] = KSone(data,func)
2 % PARÁMETROS DE ENTRADA: un vector data y una función anónima func de una única variable que
3 % representa la distribución cumulativa teórica que suponemos tiene la población
4 % de la que se han tomado los datos.
5 % PARÁMETROS DE SALIDA: la máxima separación d entre los datos contenidos en el
6 % vector data y la probabilidad prob de que el valor de d se seba al azar.
7 % Un valor pequeño de prob significa que la distribución cumulativa de los
8 % datos en data es muy probablemente diferente de func (es decir, que las diferencias
9 % observadas no se deben al azar)
10 % Esta función está tomada del capítulo 14.3 del Numerical Recipes en C
11 - data=data(:); % Convierte a data en un vector columna sea cual sea su forma inicial
12 - n = length(data);
13 - fo=0.0;
14 - data=sort(data);
15 - en=n;
16 - d=0.0;
17 - for j=1:1:n % Loop over the sorted data points.
18 -     fn=j/en; % Data's c.d.f. after this step.
19 -     ff=func(data(j)); % Compare to the user-supplied function.
20 -     dt=max(abs(fo-ff),abs(fn-ff)); % Maximum distance.
21 -     if (dt > d)
22 -         d=dt;
23 -     end
24 -     fo=fn;
25 - end
26 - en=sqrt(en);
27 - prob=probks((en+0.12+0.11/en)*d); % Compute significance.
28 - end
29
30 function [sum] = probks(alam)
31 % Kolmogorov-Smirnov probability function.
32 - EPS1=0.001;
33 - EPS2=1e-8;
34 - fac=2.0;
35 - sum=0.0;
```

¿Tienen dos muestras la misma distribución?

Supongamos que hemos tomado una muestra de N_1 valores de x y hemos representado su distribución acumulativa $F_{\text{exp1}}(x)$.

Ahora tomamos otra muestra de N_2 valores de x y representamos su distribución acumulativa $F_{\text{exp2}}(x)$

En general $F_{\text{exp1}}(1) \neq F_{\text{exp2}}(x)$

Se nos plantea la pregunta:

¿Es la diferencia fruto del azar, o es que ha cambiado la población entre las dos medidas?

KS-test para dos muestras

1) Calculamos:

$$D_o = \max |F_{\text{exp1}}(x) - F_{\text{exp2}}(x)|$$

Puede demostrarse que la probabilidad de haber obtenido un valor de la distancia $D > D_o$ viene dada por:

$$P(D > D_o) = 2 \sum_{j=1}^{\infty} (-1)^{j-1} \exp(-2j^2 z^2)$$

$$z \approx \left[\sqrt{N_{\text{eff}}} + 0.12 + \frac{0.11}{\sqrt{N_{\text{eff}}}} \right] D_o \quad N_{\text{eff}} = \frac{N_1 N_2}{N_1 + N_2}$$

Código

```
Editor - F:\docencia\clases\materiales recurrentes\MNS\ analisis\KStwo.m
File Edit Text Go Cell Tools Debug Desktop Window Help
[Icons] Stack: Base
- 1.0 + ÷ 1.1 × % %
1 function [prob,d] = KStwo(data1,data2)
2 % PARÁMETROS DE ENTRADA: dos vectores data1 y data2
3 % con los valores de dos muestras diferentes que queremos comparar para ver
4 % si tienen la misma distribución acumulativa
5 % PARÁMETROS DE SALIDA: la máxima separación d entre los datos contenidos en el
6 % vector data y la probabilidad prob de que el valor de d se seba al azar.
7 % Un valor pequeño de prob significa que la distribución acumulativa de los
8 % datos en data es muy probablemente diferente de func (es decir, que las diferencias
9 % observadas no se deben al azar)
10 % Esta función está tomada del capítulo 14.3 del Numerical Recipes en C
11 - data1=data1(:); % Convierte a data1 en un vector columna sea cual sea su forma inicial
12 - data2=data2(:);
13 - n1=length(data1);
14 - n2=length(data2);
15 - data1=sort(data1);
16 - data2=sort(data2);
17 - d=0.0;
18 - j1=1;
19 - j2=1;
20 - fn1=0.0;
21 - fn2=0.0;
22 - while (j1 <= n1 && j2 <= n2)
23 -     d1=data1(j1);
24 -     d2=data2(j2);
25 -     if (d1 <= d2)
26 -         j1=j1+1;
27 -         fn1=j1/n1; % Next step is in data1.
28 -     end
29 -     if (d2 <= d1)
30 -         j2=j2+1;
31 -         fn2=j2/n2; % Next step is in data2.
32 -     end
33 -     dt = abs(fn2-fn1);
34 -     if (dt > d)
35 -         d=dt;
```

Tests usando la *Statistics Toolbox*

Si tenemos una instalación de Matlab con la Statistics Toolbox, ya tenemos programados los tests que hemos explicado en este tema:

Test de student: ***ttest(X)***, ***ttest2(X,Y)***

Test F: ***vartest2(X,Y)***

Test de Kolmogorov-Smirnoff: ***kstest(X)***, ***kstest2(X1,X2)***

Además de toda una serie de distribuciones cumulativas, diferenciales y sus inversas.

Bibliografía del tema:

- W.H. Press, S.A. Teukolsky, W.T. Vetterling, B.P. Flannery. Numerical Recipes in C / Fortran 77 / Fortran 90, Cambridge University Press (1997). Disponible de forma gratuita en: <http://www.nr.com/oldverswitcher.html>
- Milton Abramowitz and Irene A. Stegun . Handbook of mathematical functions with formulas, graphs and numerical tables. Dover Books, 1972. Disponible de forma gratuita en: http://www.nrbook.com/abramowitz_and_stegun/
- Joaquim P. Marqués de Sá, Applied statistics using SPSS, STATISTICA, MATLAB and R , Springer Verlag, (2007)