

# Introducción al análisis estadístico de datos de datos

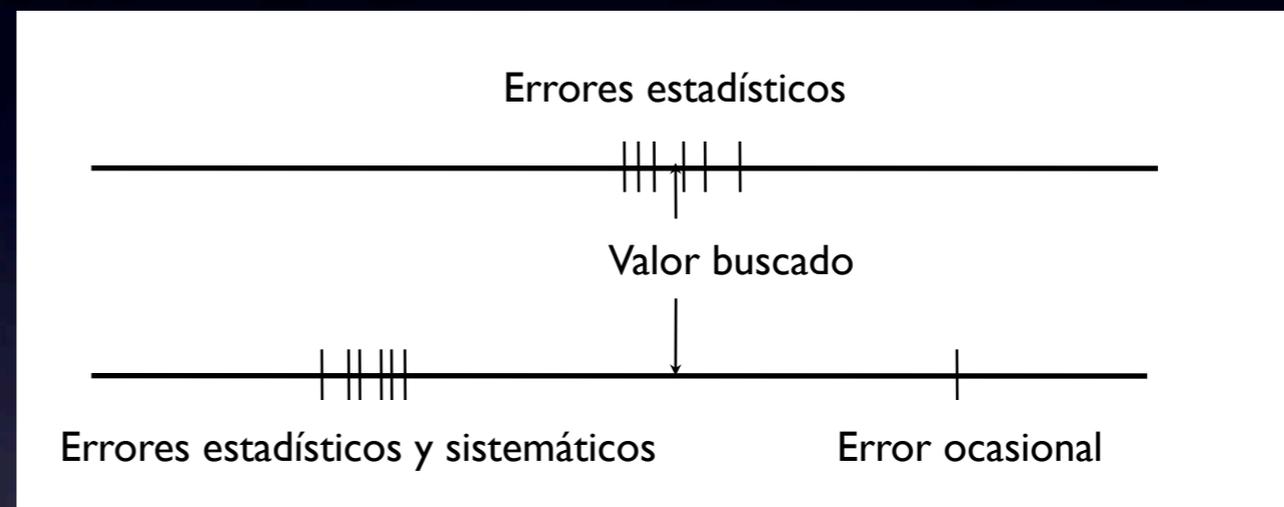
Juan A. Garzón / USC

*No hay medida sin error  
ni hortera sin transistor*

*(Refranero popular)*

# Sobre errores e incertidumbres

Toda medida está afectada de un **error**: Llamamos error a la diferencia entre el resultado de una medida y en valor buscado (o valor verdadero). Es una cantidad desconocida de la que solo podremos conocer, o hacer hipótesis, su comportamiento estadístico



- **Errores estadísticos**: se pueden hacer hipótesis sobre su comportamiento y evolución
- **Errores sistemáticos**: son muy difíciles de estimar. Se pueden acotar mediante la medida de un valor patrón

# Sobre errores e incertidumbres

La **incertidumbre** es una cantidad que permite asignar un cierto **grado de confianza** al **resultado** de una medida y determinar la **probabilidad** de que el resultado de la medida esté más o menos cerca del valor buscado.

## Ejemplo:

Se ha medido una variable  $a$  y el resultado final ha proporcionado un valor  $a=10.4$ , con una incertidumbre  $\delta a=0.4$ . Escribimos:

$$a = 10.6 \pm 0.4$$

Si suponemos que la incertidumbre tiene un comportamiento gaussiano el resultado indica que (como ya veremos) el valor que estamos midiendo:

- tiene una probabilidad del 68% de estar entre 10.2 y 11.0 (a 1 desviación típica)
- tiene una probabilidad del 95% de estar entre 9.8 y 11.4 (a 2 desviaciones típicas)
- etc.

(Interpretación inversa de la probabilidad)

- **Incertidumbres estadísticas:** Son las debidas a los errores estadísticos. Se pueden reducir aumentando el número de datos ("aumentando la estadística")
- **Incertidumbres no estadísticas:** Son las debidas a todos los errores no estadísticos (sistemáticos, ocasionales,...). Solo se pueden reducir mejorando el procedimiento de medida. Se pueden acotar midiendo un valor patrón.

En rigor, ambas incertidumbres no deben de componerse y deben de ser representadas por separado.

## Ejemplo:

$$g = 9.817 \pm (0.023)_{\text{noest}} \pm (0.015)_{\text{est}} \text{ m.s}^{-2}$$

# Propagación de incertidumbres estadísticas

## Dos variables:

Sea la variable  $x$  relacionada con las variables  $a$  y  $b$  a través de la relación analítica:

$$x = f(a,b)$$

Si se conocen las incertidumbres estadísticas  $\delta a$  y  $\delta b$ , la incertidumbre estadística de  $x$  viene dada por la relación:

$$\delta^2 x = (\partial f / \partial a)^2 \cdot \delta^2 a + (\partial f / \partial b)^2 \cdot \delta^2 b$$

La extrapolación a varias variables es inmediata

## Algunos casos particulares:

$$\begin{aligned} x = a+b \text{ o } x = a-b &\Rightarrow (\delta x)^2 = (\delta a)^2 + (\delta b)^2 \\ x = a \cdot b \text{ o } x = a/b &\Rightarrow (\delta x/x)^2 = (\delta a/a)^2 + (\delta b/b)^2 \end{aligned}$$

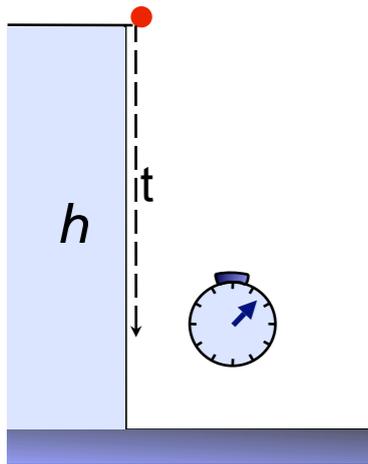
## Ejemplos :

$$x = c^2 = c \cdot c \quad \left(\frac{\partial x}{x}\right)^2 = \left(\frac{\partial c}{c}\right)^2 + \left(\frac{\partial c}{c}\right)^2 = 2 \left(\frac{\partial c}{c}\right)^2 = 2 \frac{\partial^2 c}{c^2}$$

$$x = \frac{a+b}{c^2} \quad \left(\frac{\partial x}{x}\right)^2 = \left(\frac{\partial(a+b)}{a+b}\right)^2 + \left(\frac{\partial c^2}{c^2}\right)^2 = \frac{\partial^2(a+b)}{(a+b)^2} + 2 \left(\frac{\partial c}{c}\right)^2 = \frac{\partial^2 a + \partial^2 b}{(a+b)^2} + 2 \frac{\partial^2 c}{c^2}$$

# Incertidumbres estadísticas y sistemáticas

## Ejemplo:



Se desea calcular la aceleración de la gravedad mediante un experimento de caída libre. Para ello se deja caer un cuerpo desde una altura  $h=2.100m$  y se mide varias veces el tiempo de llegada al suelo obteniéndose los siguientes resultados, en segundos:

0.654, 0.658, 0.655, 0.652, 0.654, 0.654, 0.655, 0.652, 0.653 y 0.654

El valor medio de las medidas es:  $\langle t \rangle = 1/n \cdot \sum t_i = 0.6541$

La dispersión de las medidas es:  $\sigma_{n-1} = \sqrt{1/(n-1) \cdot (t_i - \langle t \rangle)^2} = 0.0017$

La dispersión o incertidumbre en  $\langle t \rangle$   $\sigma_t = 1/\sqrt{n} \cdot \sigma_{n-1} = 0.0005$

## Fórmulas:

$$h = \frac{1}{2} g \cdot t^2 \quad \Rightarrow \quad g = \frac{2h}{t^2}$$

$$g = 2 \cdot 2.100m / 0.6541^2 s^2 = 9.817 m \cdot s^{-2}$$

## Incertidumbres:

$$\left(\frac{\partial g}{g}\right)^2 = \left(\frac{\partial h}{h}\right)^2 + 2\left(\frac{\partial t}{t}\right)^2$$

La incertidumbre de la altura  $h$  suponemos que es "no estadística". No se dice nada de ella así que suponemos que solo se ha medido una vez con una precisión de  $1mm$  (última cifra significativa). En este caso, la resolución está dada por:

$$\partial h = 0.001m / \sqrt{12} = 0.0003m$$

y:

$$\partial g_{NoEstad} = g \cdot dh / h = 9.8 \cdot 0.0003 / 2.1 = 0.0014$$

$$\partial g_{Estad} = g \cdot \sqrt{2} dt / t = 9.8 \cdot 1.4 \cdot 0.0005 / 0.65 = 0.011$$

Finalmente :

$$g = 9.817 (\pm 0.001)_{NoEstad} (\pm 0.011)_{Estad} m \cdot s^{-2}$$

# Sobre errores e incertidumbres

En rigor, una probabilidad (o una densidad de probabilidad) está asociada a la magnitud que se mide  $a_{med}$ , no a la que se busca  $a_{busc}$ . Cuando se escribe (suponiendo incertidumbres gaussianas) que el resultado de una medida ha sido:

$$a_{med} = 100 \pm 1$$

Su significado es

Si  $a_{busc}$  fuese mayor o igual que 99,  $a_{med}$  estaría a menos de una desviación típica de  $a_{busc}$   
Si  $a_{busc}$  fuese menor o igual que 101,  $a_{med}$  estaría a menos de una desviación típica de  $a_{busc}$

Es decir: 99 y 101 definen un intervalo tal que:

Si  $a_{busc}$  estuviese dentro de dicho intervalo, la probabilidad de haber observado nuestra medida sería mayor del 68%.

De idéntica forma se interpretarían el intervalo (98,102), asociado a un probabilidad del 95.4%, ... etc.

**Interpretación inversa de la probabilidad:**

En un ejemplo como el anterior, es habitual decir:

El valor buscado  $a$  tiene una probabilidad del 68% de encontrarse en el intervalo (99,101)

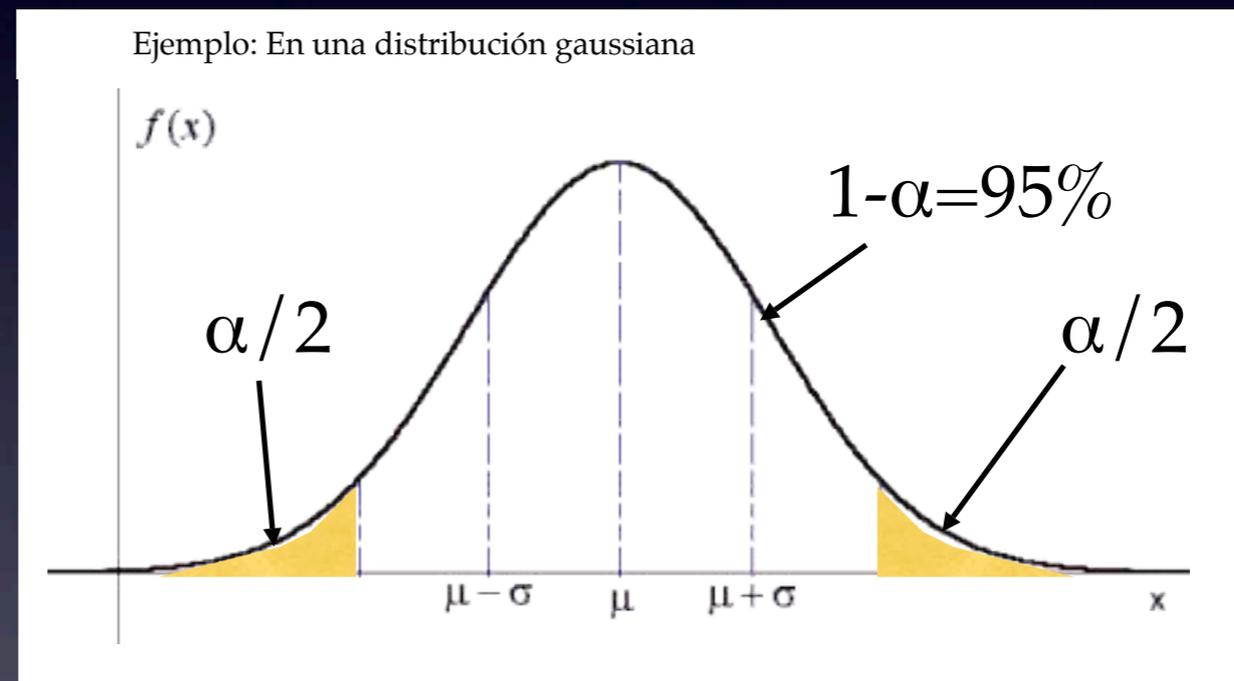
Sin embargo, dicha así, la expresión es incorrecta, puesto que el valor buscado está o no está en dicho intervalo y no se le puede asociar una probabilidad

# Sobre intervalos e índices de confianza

Alrededor del valor medio de una distribución se pueden definir un **intervalo de confianza** relacionado con la probabilidad de que una medida esté o no dentro de dicho intervalo (que no tiene por que estar centrado en el valor medio).

A la probabilidad asociada a dicho intervalo se le denomina **nivel de confianza** o  $(1-\alpha)$ .

A  $\alpha$  se la denomina, a veces, "significación estadística".



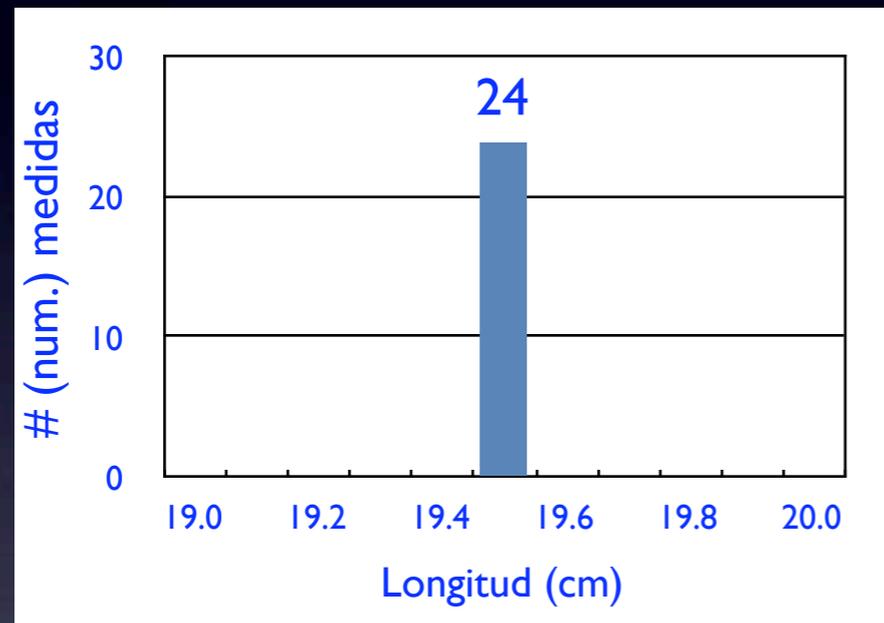
Así, en una distribución gaussiana, al intervalo centrado en su valor medio con un anchura de  $4\sigma$  le corresponde un nivel de confianza del 95%.

En el caso de distribuciones no simétricas (como la distribución de Poisson) dado un cierto nivel de confianza, el intervalo de confianza puede ser elegido no simétrico. En ese caso, se puede elegir un intervalo "centrado" o bien aquel que sea "menor".

# La medida

Ej: Medida de la longitud de una varilla con una regla  
Se dan varias posibilidades:

a: La medida proporciona siempre el mismo valor

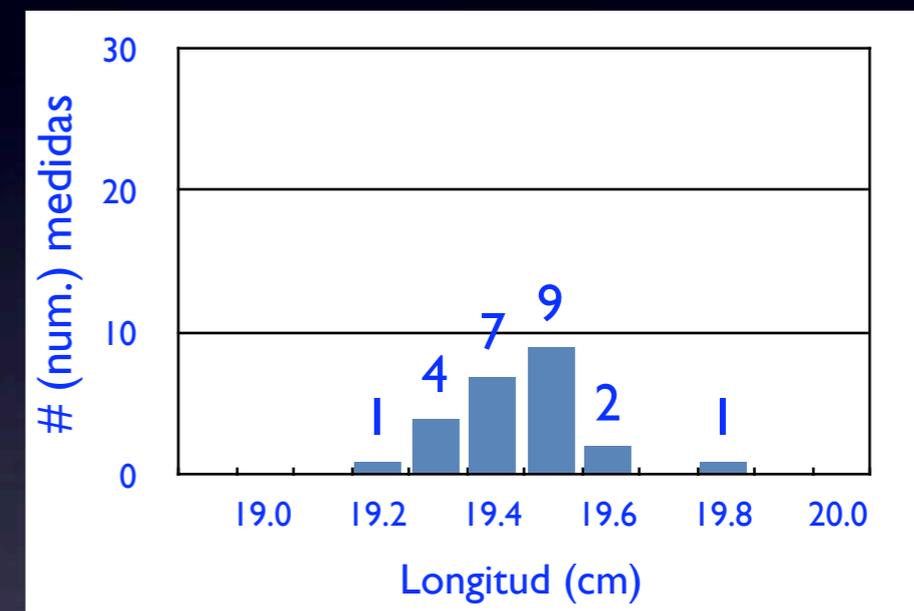


Aparentemente es el caso mas favorable

Sin embargo:

- La repetición de la medida no aporta ventajas
- No se puede obtener mejor resultado que el proporcionado por la precisión del aparato de medida: el intervalo mínimo de medida  
(Ej:  $L=19.5$ )

b: La medida proporciona distintos valores



Aparentemente es el caso mas desfavorable

Sin embargo (bajo ciertas condiciones):

- La repetición de la medida permite mejorar el conocimiento de la magnitud medida
- La calidad (incertidumbre) de la medida se puede mejorar tanto como se quiera aumentando el número de medidas

# Tablas, Histogramas y distribuciones

Cuando se dispone de una serie de medidas de una misma magnitud el primer paso es agruparlas en una tabla y representarlas en forma de histograma o distribución.

Ejemplo:

Sea un conjunto de 24 medidas de la longitud de una varilla.

Existen diversas formas de agrupar los datos:

Inicio de intervalo(cm)	Anchura del intervalo(cm)	Num. de entradas	Num. Entradas/ Anchura Intervalo	Num. Entradas/ Anchura.NumSucesos
19,0	0,1	0	0,0	0,00
19,1	0,1	0	0,0	0,00
19,2	0,1	1	10,0	0,40
19,3	0,1	4	40,0	1,70
19,4	0,1	7	70,0	2,90
19,5	0,1	9	90,0	3,80
19,6	0,1	2	20,0	0,80
19,7	0,1	0	0,0	0,00
19,8	0,1	1	10,0	0,40
19,9	0,1	0	0,0	0,00
	Total:	24	240,0	10,00

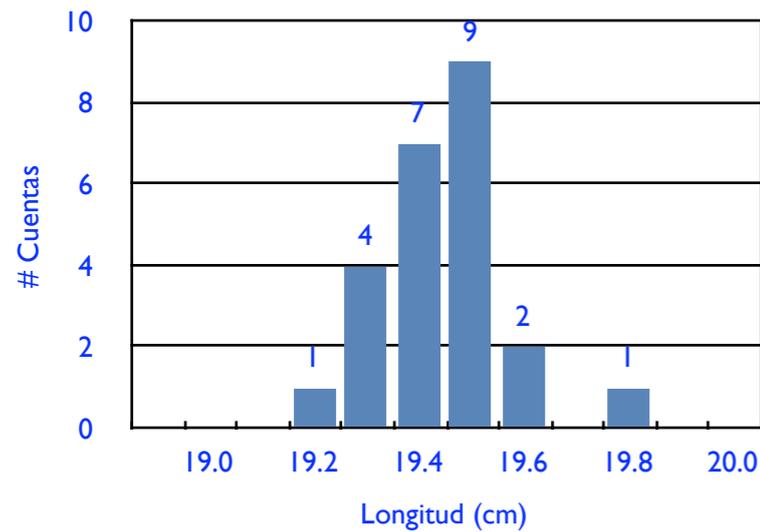
↓  
Histograma

↘  
Distribuciones

# Tablas, Histogramas y Distribuciones

## Histograma:

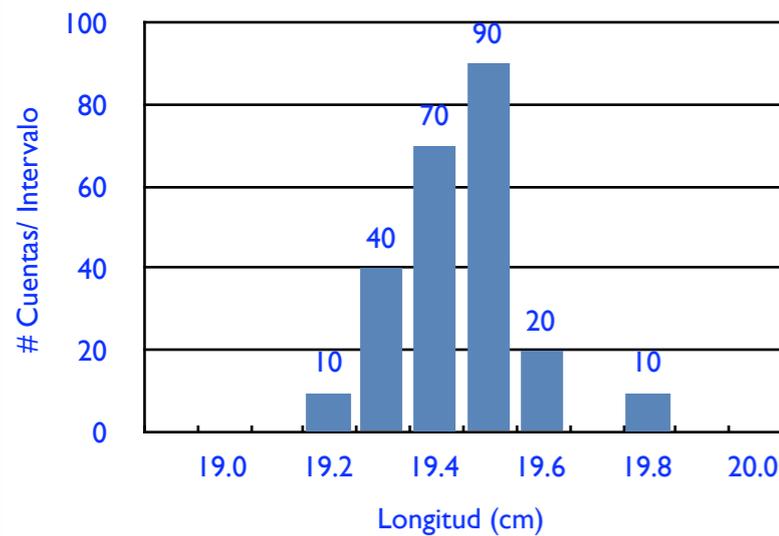
Suma de altura de las barras = Nsucs



## Distribución:

Integral (suma de superficie de las barras) = Nsucs

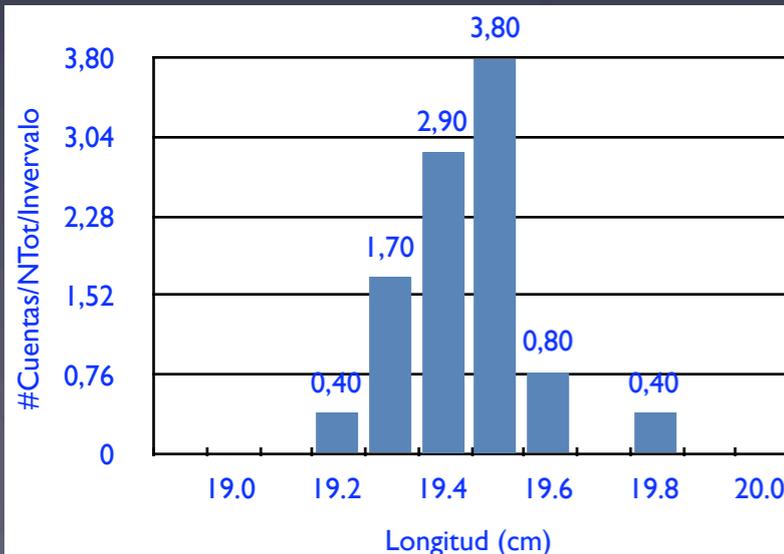
(Esta será la distribución que deba de compararse con una distribución matemática no normalizada)



## Distribución normalizada:

Integral (suma de superficie de las barras) = 1

(Esta será la distribución que deba de compararse con una distribución matemática normalizada)



# Medidas características de una distribución

Sea un conjunto de N medidas:  $x_1, x_2, x_3 \dots x_N$

## Medidas de centralización:

- **Valor medio o media aritmética:**

$$\bar{x} = \sum_{i=1}^N x_i / N$$

Si los datos están agrupados en k intervalos, cada uno centrado en  $m_j$  y con  $n_j$  entradas:

$$\bar{x} = \sum_{j=1}^k (n_j / N) m_j \quad \text{o} \quad \bar{x} = \sum_{j=1}^k p_j m_j$$

siendo  $n_j / N = p_j$  la probabilidad (frecuentista) de obtener una medida en el intervalo j  
(nota: ambas definiciones de valor medio pueden diferir ligeramente)

- **Mediana,  $x_M$ :** Es aquel valor tal que, ordenados en magnitud los datos, hay un mismo número de ellos por encima y por debajo de dicho valor
- **Moda,  $x_m$ :** Es el valor mas frecuente

# Medidas características de una distribución

Sea un conjunto de N medidas:  $x_1, x_2, x_3 \dots x_N$

## Medidas de dispersión:

- **Desviación típica:**

$$s = \sqrt{\sum_{i=1}^N (x_i - \bar{x})^2 / N}$$

o 
$$s_{N-1} = \sqrt{\sum_{i=1}^N (x_i - \bar{x})^2 / (N-1)}$$

Si los datos están agrupados en intervalos o clases,

$$s = \sqrt{\sum_{j=1}^k (x_j - \bar{x})^2 \cdot p_j}$$

- **Varianza:** Es el cuadrado de la desviación típica:

$$V = s^2$$

## Otras medidas:

- **Residuo:**

$$r_i = x_i - \bar{x}$$

- **Momento de orden k:**

$$\mu_k = 1 / N \sum_{j=1}^N (x_j - \bar{x})^k$$

Nota: Los momentos son también medidas de dispersión: La varianza es el momento de orden 2:  $V = \mu_2$

# Medidas características de una distribución

## Algunos comentarios:

Las medidas características permiten reducir un número muy grande de datos a unos pocos valores: el valor de centralidad, el de dispersión.... Si se conoce la distribución matemática de los datos, la pérdida de información es muy pequeña e, incluso, nula.

- El **valor medio** es quizás la medida de centralidad mas utilizada.
  - Es fácil de calcular y de programar en las calculadoras de bolsillo por métodos recurrentes. Utiliza la información de todos los datos.
  - Sin embargo: Es muy sensible a algún dato anómalo.

No verifica, en general, que:  $\overline{x^2} = \bar{x}^2$

- La **mediana** solo tiene en cuenta el orden de los datos, y no su magnitud, por lo que no se altera mucho si alguna observación tiene un gran error.

Sí verifica, por ejemplo, que  $x^2_M = x_M^2$

- El uso de los ordenadores han facilitado su cálculo por lo que su uso se está extendiendo
- La **moda** es útil en aquellos casos en que los datos siguen una distribución matemática no normalizable  
(P. Ej.\_ Distribución de Breit-Wigner o de Cauchy)
- Los tres valores son aproximadamente iguales si los datos tienen una distribución casi simétrica

Las medidas características dependen del tipo de distribución que siguen los datos y apenas cambian cuando el número de medidas se hace muy grande

# Medidas características de una distribución

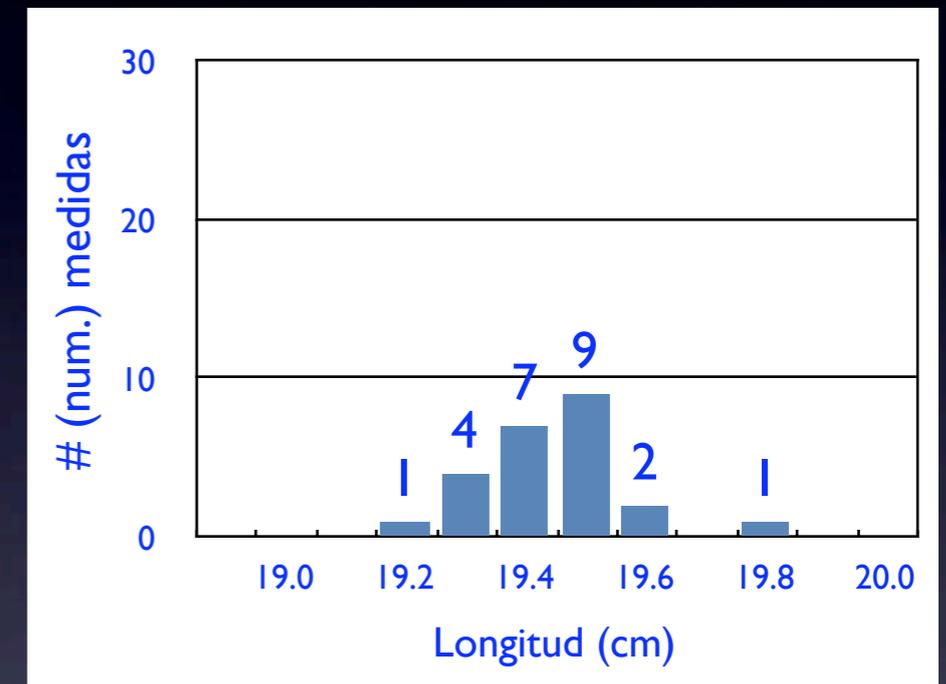
## Ejemplo:

Sea el caso de un conjunto de datos que se han ido agrupando en un histograma según se tomaban (no se dispone de los valores individuales)

Utilizamos las fórmulas para datos agrupados:

Valor medio:  $\bar{x} = 19.45$   
Mediana:  $x_M = 19.45 = (19.4 + 19.5)/2$   
Moda:  $x_m = 19.5$

Varianza:  $V = 0.15$   
Desviación típica:  $s = 0.12$



### Comentarios:

- Obsérvese que se verifica la **desigualdad de Tchebychev**, que dice:

El número de medidas entre el valor medio  $\bar{x}$  y  $\pm k$  veces la desviación típica  $s$  es como mínimo:

$$100\left(1 - \frac{1}{k^2}\right)\%$$

Ej: Para  $k=2$ , entre 19.18 y 19.66, hay 21 medidas ( $> 18 = 24 \cdot [1 - 0.25]$ )

# Análisis estadístico de datos

El análisis de datos permite extraer de una medida o de un conjunto de medidas la información que requiere el observador. Tras una toma de datos, suele darse algunos de los casos siguientes:

- Los datos se distribuyen siguiendo alguna cierta distribución matemática conocida

El problema consiste en determinar el o los parámetros característicos de la distribución y su incertidumbre: un número que permita estimar (en ausencia de errores sistemáticos graves) cuán probable es que los valores medidos se encuentren a una cierta distancia de los valores buscados: Ej. **Distribuciones binomial, Poisson, gaussiana, Chi2...**

- Los datos se distribuyen siguiendo una distribución física para la que se puede disponer o no de un modelo físico

El problema consiste en determinar la mejor forma funcional de la distribución, y los parámetros que la definen, con sus incertidumbres, asignando algún índice de calidad (o de confianza) a las diversas formas funcionales propuestas. Ej. **Distribución de masa invariante de los productos de la desintegración de una partícula inestable, los niveles de energía en un espectro....**

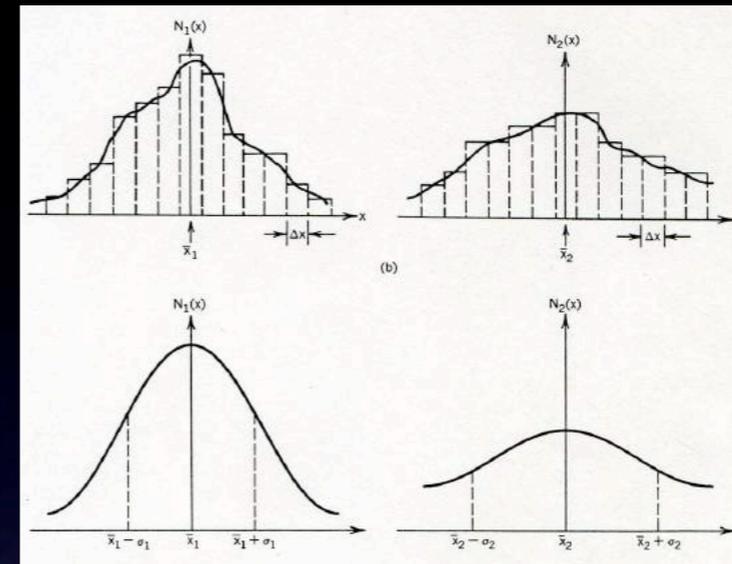
- El resultado ha sido negativo y no se tienen datos

En este caso, si existe algún modelo previo, el problema consiste en estimar cotas a los parámetros del modelo, con sus incertidumbres: Ej. **Probabilidad de desintegración del protón**

# Distribuciones matemáticas y distribuciones físicas

## Distribuciones matemáticas:

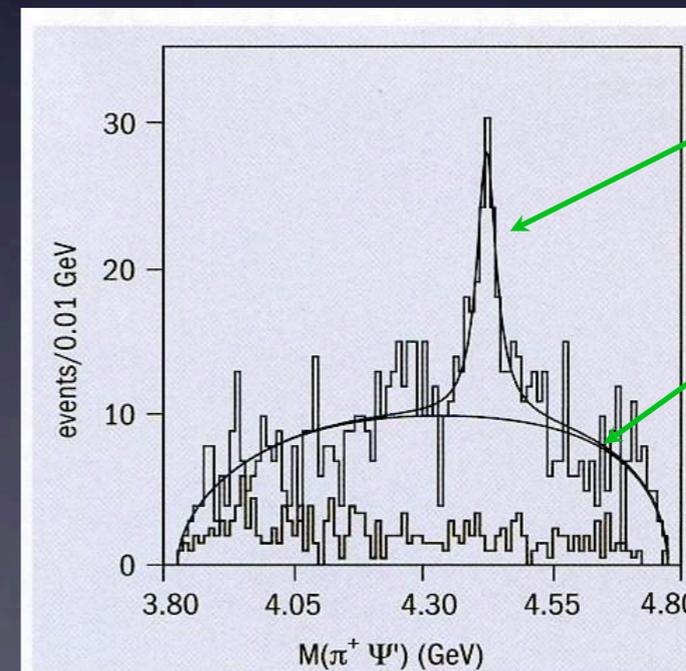
Los datos se distribuyen siguiendo alguna distribución matemática conocida:



Ejemplos de la medida de una misma variable física con un comportamiento aproximadamente gaussiano. La medida de la izquierda tiene mejor resolución (menor dispersión) que la medida de derecha)

## Distribuciones físicas:

- Los datos se distribuyen siguiendo una distribución física para la que se puede disponer o no, totalmente o parcialmene, de un modelo físico:



Resonancia:  
Distribución de Breit-Wigner

Fondo polinómico

This diagram shows the distributions of the measured mass for combinations of a detected  $\pi$  meson and a  $\Psi'$  meson. The newly discovered Z(4430) is visible as the peak at 4430 MeV.

# El logbook (o cuaderno de bitácora)

El logbook es un cuaderno donde se escriben los datos de un experimento, se anotan las incidencias, comentarios...

Normalmente consta la siguiente información:

- Fecha y hora de la toma de datos y nombre de los investigadores:
- Disposición geométrica del experimento y cualquier consideración que se considere útil para el posterior análisis.
- Lista de datos que toman
- Análisis previo e inmediato de los datos.

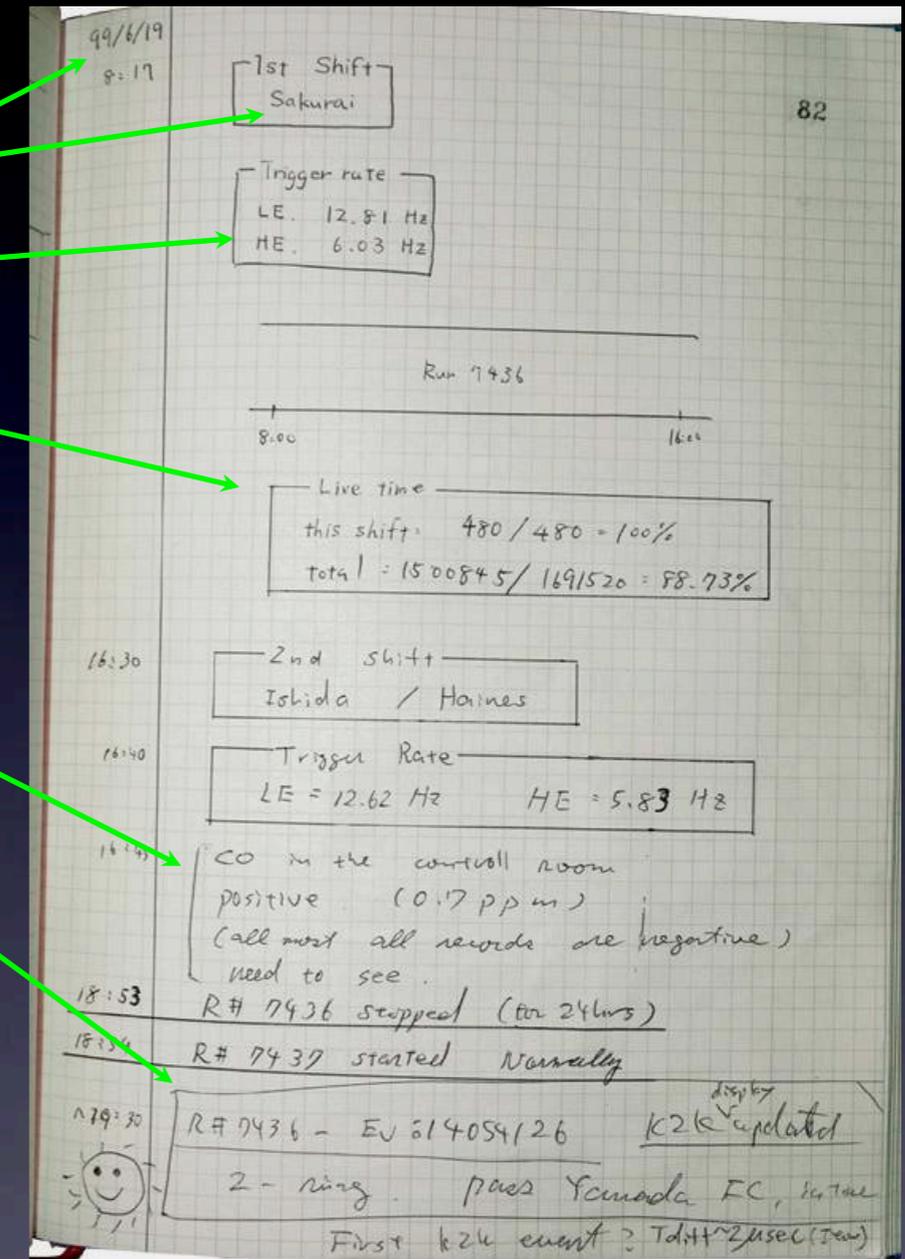
Este aspecto es muy importante: El análisis previo muestra que la toma de datos está siendo correcta y que puede continuarse. En el caso en que se detecte algún problema, es el momento de modificar la configuración del experimento, tomar datos durante mas tiempo, modificar los intervalos de medida...

Consejos:

En el logbook debe de escribirse con la mayor claridad posible para que cualquier investigador (incluido el que hace las anotaciones; olvidar es muy fácil) pueda entender con posterioridad el trabajo realizado

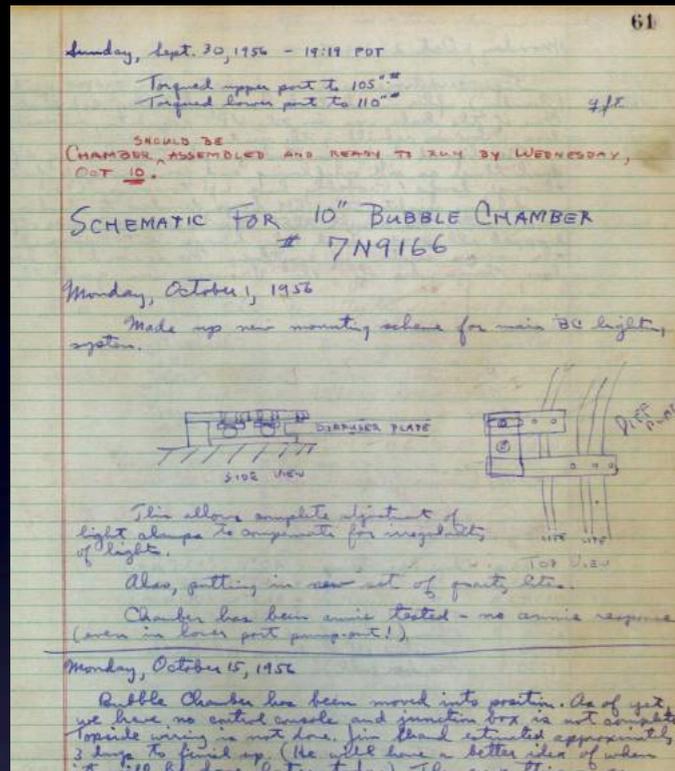
Es preferible ser redundante a ser parco en los comentarios. Toda información (fuente utilizada, energía de la radiación, fecha de calibración, tensión del detector, umbrales en la electrónica, etc.) puede ser útil con posterioridad

En el caso de contajes, ajustar los tiempos de medida de acuerdo con el ritmo de cuentas. Si se repiten medidas y los resultados entre ellas son compatibles sumar los resultados para mejorar la incertidumbre



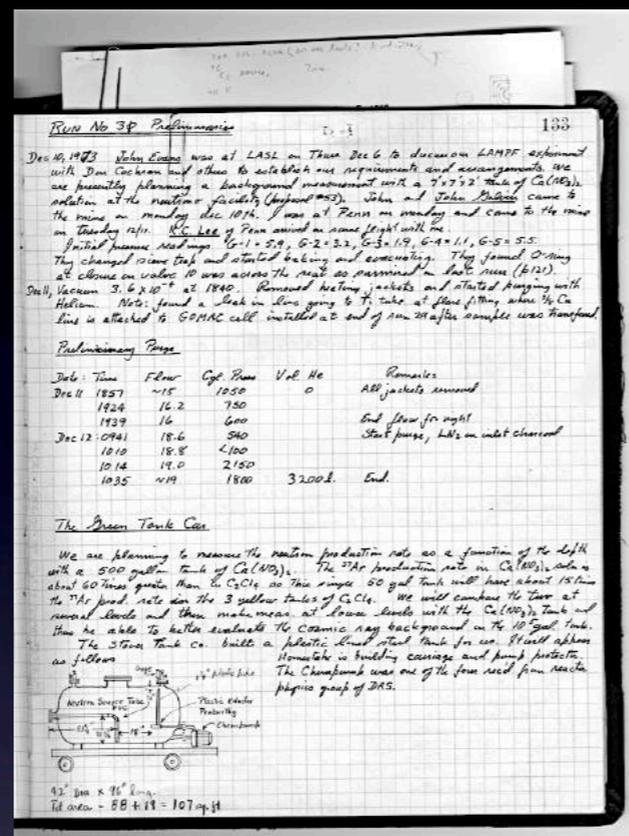
Página del logbook de un experimento en el detector de neutrinos subterráneo de Kamiokande en el momento de detectar un neutrino proveniente del acelerador KEK, situado a más de 250km de distancia

# Algunos logbooks famosos



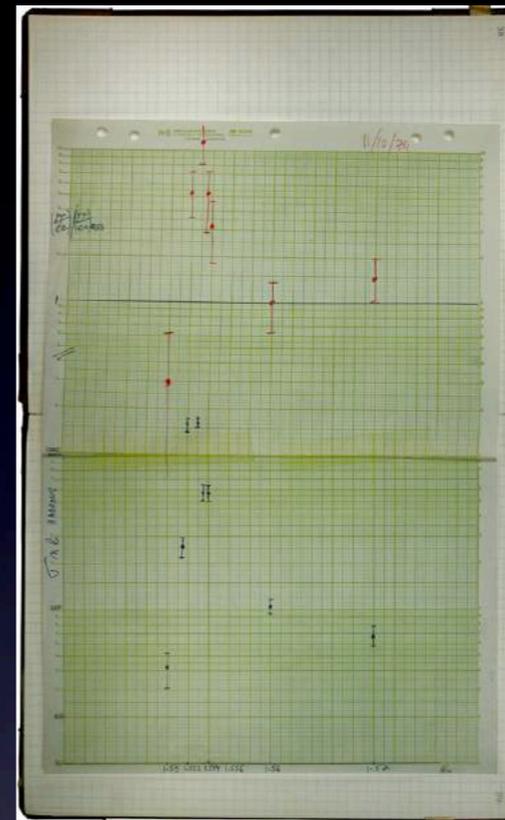
1956: Luis Alvarez

Cuaderno describiendo el diseño de un cámara de burbujas de su invención.



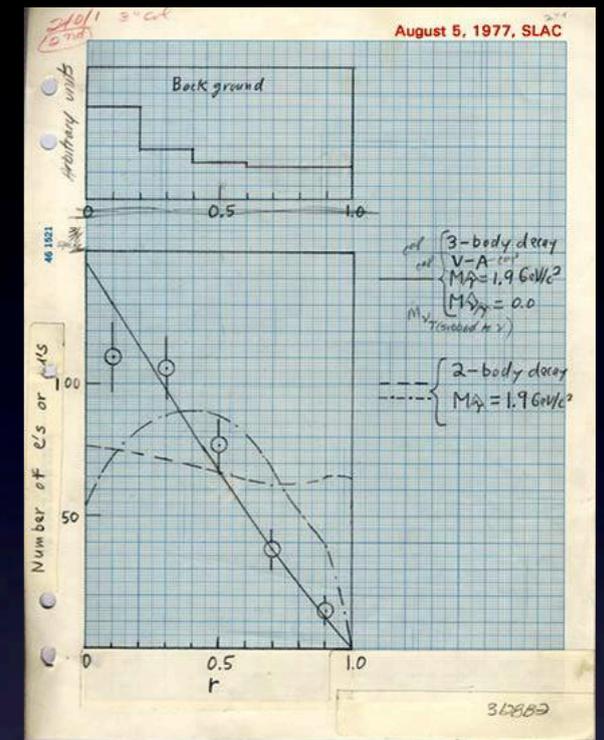
1973: Raymond Davis

Cuaderno con comentarios sobre un tanque con nitrato de calcio para estudiar el fondo de neutrones a distintas profundidades de la mina.



1974: Burton Richter

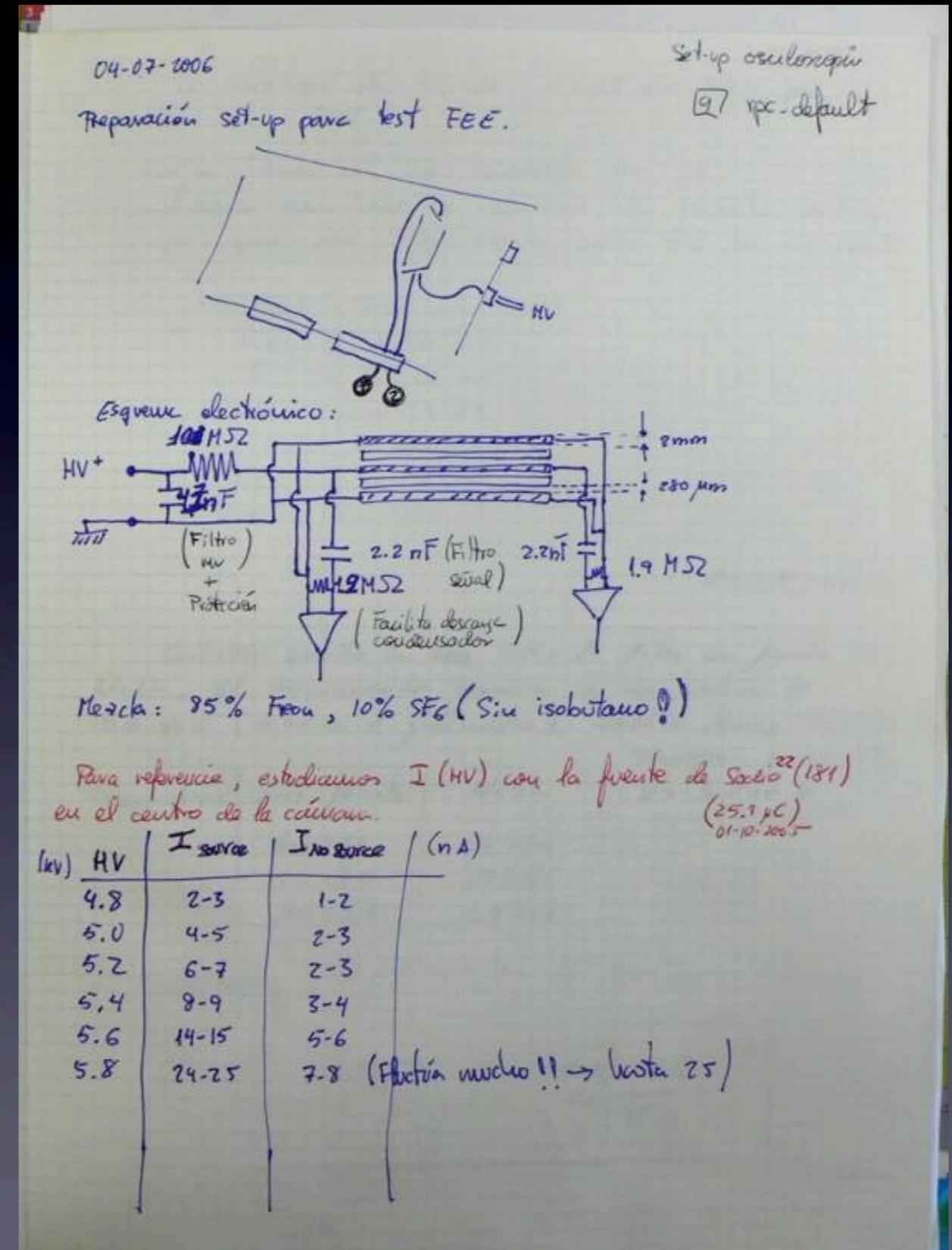
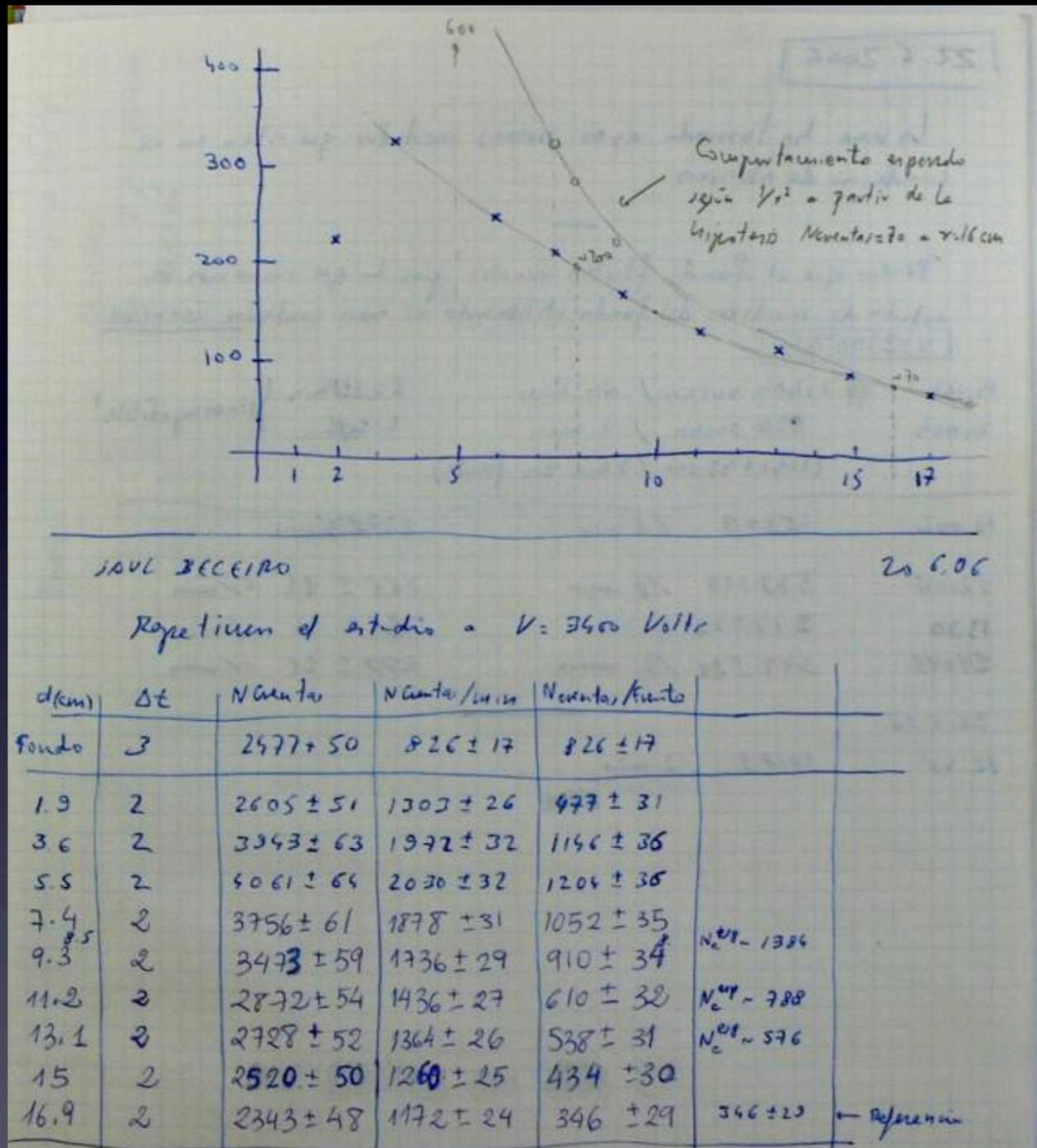
Descubrimiento de la partícula  $J/\psi$   
Hoja mostrando el número de hadrones producidos en la aniquilación  $e^+e^-$  en función de la energía en el centro de masas



1977: Martin Perl

Descubrimiento del lepton tau.  
Cuaderno mostrando los datos experimentales obtenidos y las curvas correspondientes a diversos decaimientos en dos o tres partículas.

# Mas logbooks menos famosos



# Probabilidad y Densidad de Probabilidad

## Probabilidad:

Sea un conjunto discreto de posibles medidas exclusivas de una variable  $a$ :  
 $A = \{a_1, a_2, \dots, a_k\}$ .

La función real  $P$  es una probabilidad si se verifica:

a)  $P(a_i) \geq 0 \quad \forall i$

b)  $P(a_i \text{ o } a_j) = P(a_i) + P(a_j)$

c)  $\sum_{i=1}^k P(a_i) = 1$

(Axiomas de Kolmogorov)

Ejemplo: ¿Cuál es la probabilidad de obtener cara (c) o cruz (f) lanzando una moneda al aire?

1. Por b) y c)  $\Rightarrow P(c) + P(f) = 1$

2. Por hipótesis de simetría:  $P(c) = P(f)$

1. y 2.  $\Rightarrow P(c) = P(f) = 0.5$

*Problema: ¿Qué ocurre si la moneda no es simétrica y no es cierta la hipótesis 2.?*

## Probabilidad (Definición frecuentista)

Dado un conjunto discreto de  $N$  medidas de una variable  $a$ :  $A = \{a_1, a_2, \dots, a_k\}$ , se define:

$$P(a_i) = n_i/N$$

siendo  $n_i$  el número de observaciones del resultado  $a = a_i$ .

Ejemplo: ¿Cuál es la probabilidad de obtener cara (c) o cruz (f), lanzando una moneda al aire, si ha lanzado 100 veces y se han obtenido 54 veces cara y 46 veces cruz?

$$P(c) = 54/100 = 0.54 \quad \text{y} \quad P(f) = 46/100 = 0.46$$

*Problema: Para obtener un resultado mas preciso hay que realizar mas medidas (y ello no siempre es posible). ¿Qué ocurre si estamos sesgando los lanzamientos y el recuento no es correcto?*

# Probabilidad y Densidad de Probabilidad

## Función densidad de probabilidad: $f(x)$

Sea  $A=(x_i, x_f)$  un intervalo continuo de posibles medidas exclusivas de una variable  $x$ . Una función real  $f(x)$  es una función de densidad de probabilidad si verifica:

$$a) f(x) \geq 0 \quad \forall x \in (x_i, x_f)$$

$$b) \int_{x_i}^{x_f} f(x) dx = 1$$

Por extensión se puede definir la probabilidad de un subconjunto  $B \subset A$   $P(B)$ , como:

$$P(B) = \int_B f(x) dx$$

## Función de distribución o Función acumulativa: $F(x)$

Se define como:

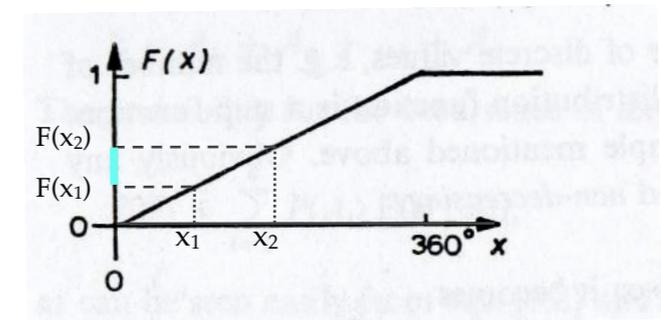
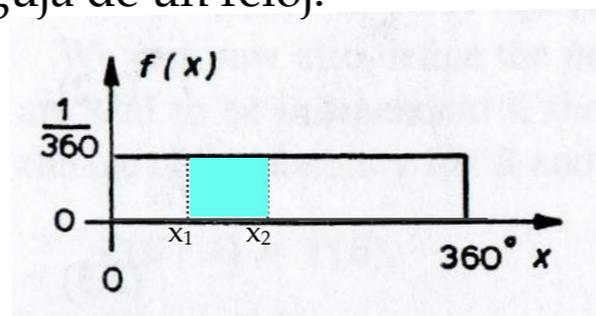
$$F(x) = \int_{x_i}^x f(x) dx$$

Por definición:

$$P(x_1, x_2) = \int_{x_1}^{x_2} f(x) dx$$

Ejemplo:

Función de densidad  $f(x)$  y función de distribución  $F(x)$  de la posición de una aguja de un reloj.



# Parámetros característicos de distribuciones matemáticas

Al igual que para las distribuciones experimentales de datos, a las distribuciones matemáticas de probabilidad (discretas y continuas) se le pueden asignar las medidas de centralidad y de dispersión correspondientes:

	Distribuciones discretas	Distribuciones continuas
Valor esperado o media: $E(x)$	$\mu = \sum_{i=1}^N x_i P(x_i)$	$\mu = \int_{-\infty}^{\infty} x f(x) dx$
Varianza: $V=E[(x-\mu)^2]$	$V = \sum_{i=1}^N (x_i - \mu)^2 P(x_i)$	$V = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$
Desviación típica: $\sigma$	$\sigma = \sqrt{V}$	$\sigma = \sqrt{V}$

# Medidas características de distribuciones matemáticas

Los parámetros característicos de la distribución de unos datos son “estimadores” de los parámetros característicos de la correspondiente distribución matemática.

Así:

$\bar{x}$  es un estimador de  $\mu$

$V_{\text{dat.}}$  es un estimador de  $V_{\text{mat.}}$

En el caso de datos con distribución gaussiana, el valor medio  $\bar{x}$  se distribuye alrededor del valor esperado  $\mu$  según una distribución gaussiana con dispersión:

$$\delta x = \frac{1}{\sqrt{N}} s_{N-1}$$

siendo  $N$  el número de medidas.

Obsérvese que  $s_{N-1}$  es una característica de la distribución experimental y permanece aproximadamente constante cuando  $N$  se hace muy grande.

$\delta x$  decrece como  $1/\sqrt{N}$  y se hace más pequeño según aumenta el número de medidas. Por eso, es beneficioso aumentar el número de medidas (aunque llega un momento en que no es práctico y es mejor mejorar otros aspectos del experimento)

$s_{N-1}$  es un estimador centrado de la dispersión matemática  $\sigma_{\text{mat.}}$ . El factor  $(N-1)$  corrige el hecho de estar calculado con el valor medio de  $x$  en vez de con el valor esperado  $\mu$ .

# Probabilidad y Densidad de Probabilidad

## Probabilidad condicionada:

Sea el caso en que se miden dos variables  $a$  y  $b$  con valores posibles, no exclusivos, contenidos respectivamente en los conjuntos  $A=\{a_1, a_2 \dots a_k\}$  y  $B=\{b_1, b_2 \dots b_m\}$ .

Se define  $P(a_i|b_j)$ , o la **probabilidad de  $a_i$  condicionada a que se ha producido  $b_j$** , a aquella cantidad que verifica:

$$P(a_i \text{ y } b_j) = P(a_i|b_j) \cdot P(b_j) = P(b_j|a_i) P(a_i)$$

Si  $P(a_i|b_j) = P(a_i) \Rightarrow a_i \text{ y } b_j \text{ son independientes} \Rightarrow P(a_i \text{ y } b_j) = P(a_i) \cdot P(b_j)$

Ejemplo:

En la baraja española sean, por ejemplo, las siguientes observaciones:

$$a \equiv \text{oro} : P(a) = 10/40 = 1/4$$

$$b \equiv \text{rey} : P(b) = 4/40 = 1/10$$

$$a \text{ y } b \equiv \text{rey de oros} : P(a \text{ y } b) = 1/40$$

Según la definición:

$$P(a|b) = (1/40) / (1/10) = 1/4 \quad (\equiv \text{Probabilidad de ser oro, si es un rey})$$

$$P(b|a) = (1/40) / (1/4) = 1/10 \quad (\equiv \text{Probabilidad de ser rey, si es un oro})$$

Se verifica:

$$P(a|b) = P(a) \text{ y } P(b|a) = P(b) \Rightarrow \text{Ser rey y ser oro a la vez son observaciones independientes}$$

# Probabilidad y Densidad de Probabilidad

## Teorema de Bayes:

A partir de la definición de probabilidad condicional es fácil observar que:

$$P(a_i|b_j) = \frac{P(b_j|a_i) \cdot P(a_i)}{P(b_j)}$$

Esta relación permite calcular una probabilidad condicionada, conocidas las probabilidades individuales y la probabilidad condicional inversa

Ejemplo:

En la Facultad de Física de la Univ. de Santiago de Compostela, en la especialidad de Física de Partículas hay un 70% de chicos. En toda la Facultad hay el mismo número de chicas y de chicos y la proporción de estudiantes que eligen la especialidad de Física de Partículas es del 30%. ¿Qué probabilidad hay de que un cierto estudiante esté en la especialidad de Partículas?

$$\text{Probabilidad (F. Partículas | Chico)} = \text{Prob (Chico | FPartículas)} \cdot \text{Prob(FPartículas)} / \text{Prob (Chico)}$$

luego:

$$\text{Probabilidad (F. Partículas | Chico)} = 0.7 \cdot 0.3 / 0.5 = 0.42 = 42\%$$

# Análisis Bayesiano

El uso del teorema de Bayes se puede extender mas allá del concepto de probabilidad y aplicarlo en el **llamado test de hipótesis**. Esta técnica estadística intenta asignar una "probabilidad" a una cierta hipótesis cómo si se tratase de una variable.

Si  **$H_i$  es una hipótesis**, de entre un conjunto de posibles hipótesis válidas, y  **$a$**  es una cierta observación, o un conjunto de datos, el teorema de Bayes aplicado a una hipótesis dice:

$$P(H_i|a) = \frac{P(a|H_i) \cdot P(H_i)}{P(a)}$$

# Análisis Bayesiano

## Comentarios:

- El teorema supone aplicar una probabilidad a una hipótesis  $P(H_i)$ , lo que es algo complicado (Ej: ¿cuál es el espacio de hipótesis, como se normaliza su probabilidad...?)
- Una misma hipótesis se puede expresar de forma diferente pudiendo dar lugar a resultados diferentes. (Ej:  $\text{masa} > 0$  o  $\text{masa}^3 > 0$  son matemáticamente iguales. ¿Qué condición se aplica? )
- El análisis bayesiano permite tener en cuenta la información previa acerca de una cierta hipótesis algo que las técnicas "clásicas no hacen" (desperdiciando información) o hacen "a posteriori" (ignorando la información previa)

## Ejemplo:

Se mide una cierta magnitud  $a$ , obteniéndose el valor  $a=9\pm 3$ . En otro momento se mide nuevamente de forma independiente obteniéndose el valor  $a=10\pm 1$ . ¿Qué se puede decir del valor de  $a$ ?

Las técnicas "clásicas" consideran las dos medidas independientes. Las técnicas bayesianas tienen en cuenta la medida previa, lo que, de alguna forma, reduce el espacio de posibles los posibles resultados de  $a$  "deformando" su distribución de probabilidad y proporcionando un valor ligeramente diferente.

# Analisis Bayesiano

## Un ejemplo: La masa del neutrino del electrón:

Durante mucho tiempo diversos experimentos proporcionaban para el neutrino masas negativas. En ese caso no se puede aplicar el análisis clásico de intervalos de confianza y niveles de confianza. Es mejor proceder a un tipo de análisis alternativo conocido como análisis bayesiano.

**$\nu_e$  MASS**

Most of the data from which these limits are derived are from  $\beta^-$  decay experiments in which a  $\bar{\nu}_e$  is produced, so that they really apply to  $m_{\bar{\nu}_1}$ . Assuming  $CPT$  invariance, a limit on  $m_{\bar{\nu}_1}$  is the same as a limit on  $m_{\nu_1}$ . Results from studies of electron capture transitions, given below " $m_{\nu_1} - m_{\bar{\nu}_1}$ ", give limits on  $m_{\nu_1}$  itself. See also the Listings in the Neutrino Bounds from Astrophysics and Cosmology section.

VALUE (eV)	CL%	DOCUMENT ID	TECN	COMMENT
<b>OUR LIMIT</b>		1 PDG	94 RVUE	See the note below.
< 7.2	95	2 WEINHEIMER 93	SPEC	$^3\text{H}\beta$ decay
<11.7	95	3 HOLZSCHUH 92B	SPEC	$^3\text{H}\beta$ decay
<13.1	95	4 KAWAKAMI 91	SPEC	$^3\text{H}\beta$ decay
< 9.3	95	5 ROBERTSON 91	SPEC	$^3\text{H}\beta$ decay
• • • We do not use the following data for averages, fits, limits, etc. • • •				
< 8.0	95	6 VANDYCK 93		$m_{^3\text{H}} - m_{^3\text{He}}$
<14	95	7 DECMAN 92	CNTR	$^3\text{H}\beta$ decay
<23		AVIGNONE 90	ASTR	Supernova SN 1987A
		LOREDO 89	ASTR	SN 1987A
		ABBOTT 88	ASTR	Supernova SN 1987A
<29	95	8 KAWAKAMI 88	SPEC	Repl. by KAWAKAMI 91
17 to 40		9 SPERGEL 88	ASTR	Supernova SN 1987A
<27	95	10 BORIS 87	SPEC	$\bar{\nu}_e, ^3\text{H}\beta$ decay
<18	95	11 WILKERSON 87	SPEC	$\bar{\nu}_e, ^3\text{H}\beta$ decay
	95	12 FRITSCHI 86	SPEC	Repl. by HOLZSCHUH 92B

1 PDG 94 formal upper limit, as obtained from the  $m^2$  average in the next section, is 5.1 eV at the 95%CL. Caution is urged in interpreting this result, since the  $m^2$  average is positive with only a 3.5% probability. If the weighted average  $m^2$  were forced to zero, the limit would increase to 7.0 eV.

2 WEINHEIMER 93 is a measurement of the endpoint of the tritium  $\beta$  spectrum with a magnetic guiding field. The source is molecular tritium frozen onto an aluminum substrate.

**$\nu_e$  MASS SQUARED**

The tritium experiments actually measure mass squared. A combined limit on mass must therefore be obtained from the weighted average of the results shown here. The recent results are in strong disagreement with the earlier claims by the ITEP group [LUBIMOV 80, BORIS 87 (+ BORIS 88, erratum)] that  $m_{\nu_1}$  lies between 17 and 40 eV. The BORIS 87 result is excluded because of the controversy over the possibly large unreported systematic errors; see BERGKVIST 85B, BERGKVIST 86, SIMPSON 84, and REDONDO 89. However, the average for the new experiments given below implies only a 3.5% probability that  $m^2$  is positive. See HOLZSCHUH 92 for a review of the recent direct  $m_{\nu_1}$  measurements.

VALUE (eV <sup>2</sup> )	DOCUMENT ID	TECN	COMMENT
<b>- 54 ± 30 OUR AVERAGE</b>			
- 39 ± 34 ± 15	13 WEINHEIMER 93	SPEC	$^3\text{H}\beta$ decay
- 24 ± 48 ± 61	14 HOLZSCHUH 92B	SPEC	$^3\text{H}\beta$ decay
- 65 ± 85 ± 65	15 KAWAKAMI 91	SPEC	$\bar{\nu}_e$ , tritium
- 147 ± 68 ± 41	16 ROBERTSON 91	SPEC	$\bar{\nu}_e$ , tritium
• • • We do not use the following data for averages, fits, limits, etc. • • •			
- 72 ± 41 ± 30	17 DECMAN 92	SPEC	$^3\text{H}\beta$ decay
223 ± 244 ± 269	18 KAWAKAMI 88	SPEC	Repl. by KAWAKAMI 91
- 57 ± 453 ± 118	19 WILKERSON 87	SPEC	Repl. by ROBERTSON 91
- 11 ± 63 ± 178	20 FRITSCHI 86	SPEC	Repl. by HOLZSCHUH 92B

13 WEINHEIMER 93 is a measurement of the endpoint of the tritium  $\beta$  spectrum with a magnetic guiding field. The source is molecular tritium frozen onto an aluminum substrate.

14 HOLZSCHUH 92B source is a monolayer of tritiated hydrocarbon.

15 KAWAKAMI 91 experiment uses tritium-labeled arachidic acid.

16 ROBERTSON 91 experiment uses gaseous molecular tritium. The result is in strong disagreement with the earlier claims by the ITEP group [LUBIMOV 80, BORIS 87 (+ BORIS 88, erratum)] that  $m_{\nu_1}$  lies between 17 and 40 eV.

**PDG 1994**

**¡ $m_{\nu^2} < 0!$ !**

eV at the 95%CL. Caution is urged in interpreting this result, since the  $m^2$  average is positive with only a 3.5% probability. If the weighted average  $m^2$  were forced to zero,

**Técnicas clásicas:** La masa del neutrino (en concreto  $m_{\nu^2}$ ) se determina en el decaimiento  $\beta$  del tritio, como la masa restante de los productos de la desintegración. Debido a la dificultad de corregir los errores sistemáticos, el análisis clásico ha proporcionado valores ¡negativos! con una probabilidad de ser positiva del ¡3.5%!

**Técnicas bayesianas:** Imponen a priori la condición de que la masa sea positiva (y  $m_{\nu^2} > 0$ ). Proporcionan, por tanto, valores positivos (aunque la estimación de intervalos de probabilidad se hace mas complicada)

# Análisis bayesiano

Un ejemplo de la vida corriente:

IMDb: IMovieDatabase, la mejor base de datos sobre cine

IMDb: Su clasificación sobre las mejores películas

IMDb: Su procedimiento para valorar las películas

Rank	Rating	Title	Votes
1.	9.1	<a href="#">The Godfather</a> (1972)	262,562
2.	9.1	<a href="#">The Shawshank Redemption</a> (1994)	310,729
3.	9.0	<a href="#">The Godfather: Part II</a> (1974)	150,405
4.	8.9	<a href="#">Buono, il brutto, il cattivo, II</a> (1966)	85,417
5.	8.8	<a href="#">Pulp Fiction</a> (1994)	266,768
6.	8.8	<a href="#">Schindler's List</a> (1993)	178,988
7.	8.8	<a href="#">One Flew Over the Cuckoo's Nest</a> (1975)	133,808
8.	8.8	<a href="#">Star Wars: Episode V - The Empire Strikes Back</a> (1980)	189,117
9.	8.8	<a href="#">Casablanca</a> (1942)	117,440
10.	8.8	<a href="#">Shichinin no samurai</a> (1954)	65,664
11.	8.8	<a href="#">Star Wars</a> (1977)	228,808
12.	8.8	<a href="#">The Lord of the Rings: The Return of the King</a> (2003)	239,027
13.	8.7	<a href="#">12 Angry Men</a> (1957)	62,858
14.	8.7	<a href="#">Rear Window</a> (1954)	77,683
15.	8.7	<a href="#">Goodfellas</a> (1990)	145,323
16.	8.7	<a href="#">Cidade de Deus</a> (2002)	90,360
17.	8.7	<a href="#">Raiders of the Lost Ark</a> (1981)	162,005
18.	8.7	<a href="#">The Lord of the Rings: The Fellowship of the Ring</a> (2001)	272,608

248.	7.9	<a href="#">His Girl Friday</a> (1940)	12,992
249.	7.9	<a href="#">Olvidados, Los</a> (1950)	3,534
250.	7.9	<a href="#">Harold and Maude</a> (1971)	17,953

The formula for calculating the Top Rated 250 Titles gives a true Bayesian estimate:

$$\text{weighted rating (WR)} = (v + (v+m)) \times R + (m + (v+m)) \times C$$

where:

R = average for the movie (mean) = (Rating)

v = number of votes for the movie = (votes)

m = minimum votes required to be listed in the Top 250 (currently 1300)

C = the mean vote across the whole report (currently 6.7)

for the Top 250, only votes from regular voters are considered.

# Norma básica: Precaución y sentido común

- Precaución al tomar medidas características precisas sobre un número pequeño de intervalos: **El resultado puede depender de la elección de intervalos**

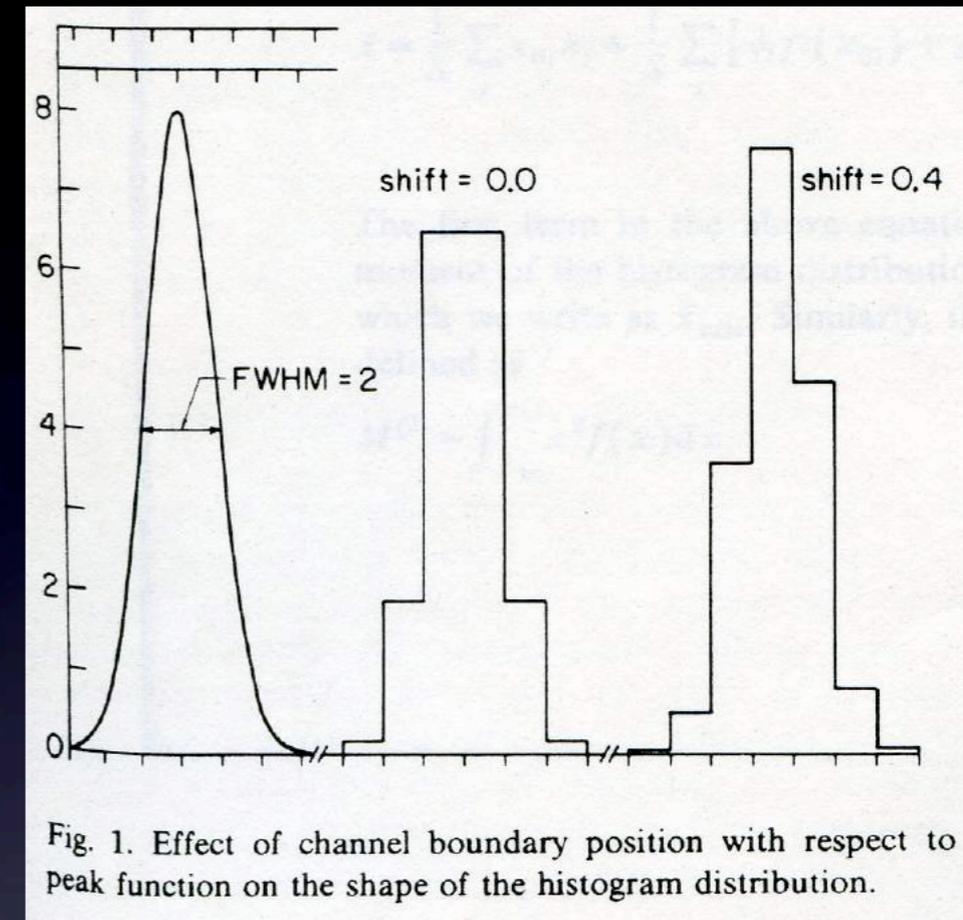
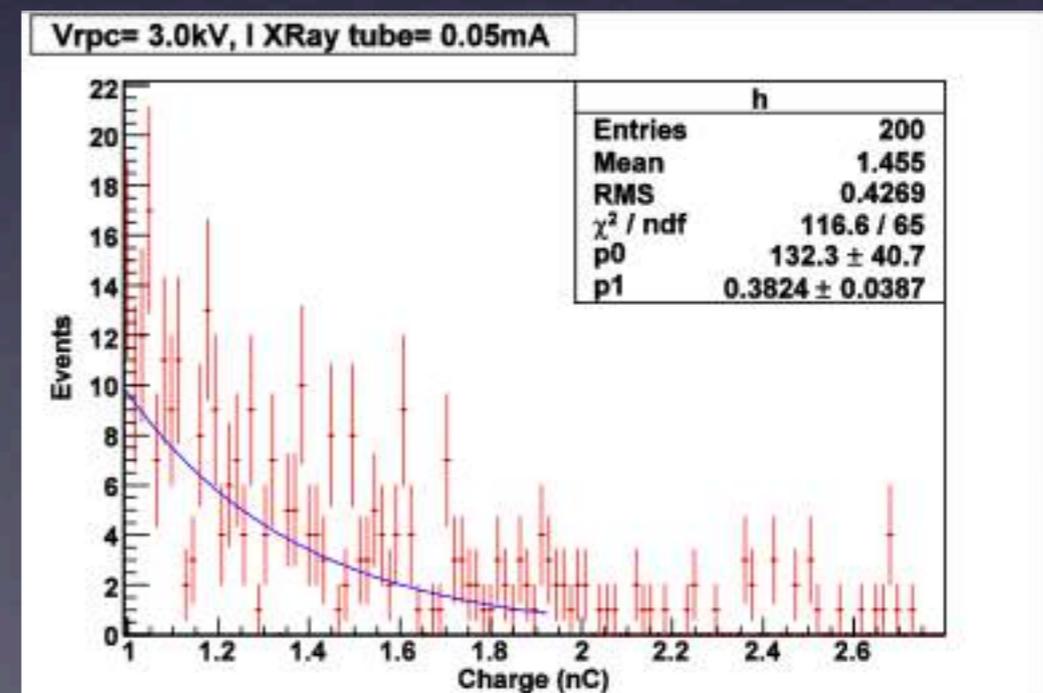


Fig. 1. Effect of channel boundary position with respect to peak function on the shape of the histogram distribution.

- Precaución al trabajar con histogramas con muchos intervalos conteniendo muy pocos sucesos. **Es mejor agrupar los intervalos y trabajar con menos intervalos aunque mas poblados.**



# Algunas distribuciones de probabilidad

# Distribución Binomial

Los recuentos de una cierta magnitud discreta se agrupan siguiendo una distribución Binomial si podemos suponer que se dan las siguientes condiciones:

1. La magnitud solo admite dos valores: Éxito o Fracaso (o Cara o Cruz, Si o No,...)
2. La probabilidad de éxito (o fracaso) es constante
3. El número de recuentos es siempre el mismo

La expresión analítica de una distribución binomial depende de **1 parámetro, p**, y es de la forma:

$$P_N(k;p) = \binom{N}{k} p^k (1-p)^{N-k}$$

Representa la probabilidad de obtener k éxitos, en N intentos, siendo la **p** la probabilidad de tener un éxito en un intento.

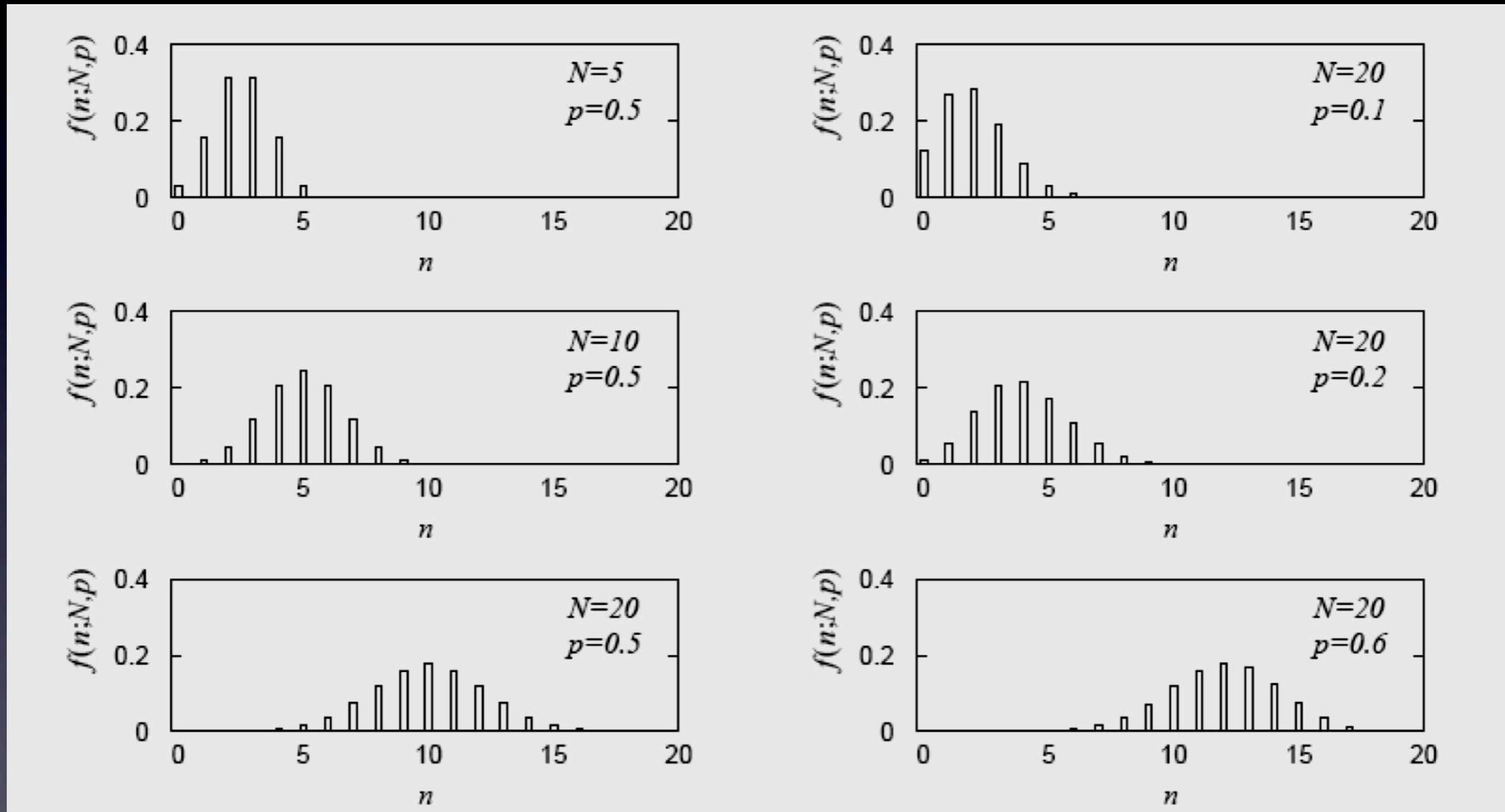
$$\begin{aligned}\mu &= N \cdot p \\ V &= N \cdot p \cdot (1-p)\end{aligned}$$

Ejemplos:

- El número de veces que salen 6 caras, con una moneda en 10 intentos
- El número de veces que sale un 2 con un dado, cada 10 intentos
- El número de veces que una partícula  $\alpha$  sale hacia atrás en el experimento de Rutherford en intervalos iguales de tiempo

# Distribución Binomial

Algunos casos de distribución binomial:



Comentario:

Todo recuento de una variable de la que solo hay dos valores posibles, y sus probabilidades son constantes en el tiempo, sigue una distribución binomial.

# Distribución de Poisson

Un conjunto de contajes (de una magnitud discreta) se agrupan siguiendo una distribución de Poisson si podemos suponer que se dan las siguientes condiciones:

1. La producción de cada cuenta es independiente y aleatoria
2. La probabilidad de cada suceso (cada cuenta) es constante
3. Los sucesos se producen a ritmo constante

La expresión analítica de una distribución de Poisson depende de **1 parámetro  $\mu$**  es de la forma:

$$P(k; \mu) = \frac{\mu^k e^{-\mu}}{k!}$$

La variable  $k$  solo toma valores enteros, aunque el parámetro  $\mu$  puede no ser entero.

$$\begin{aligned} \mu &= E(k) = \mu \\ V &= \mu \end{aligned}$$

¡El parámetro  $\mu$  es el valor esperado de las medidas!

Ejemplos:

- El número de desintegraciones radiactivas de una fuente de larga vida media por unidad de tiempo
- El número de latas de refresco por kilómetro en la cuneta de una carretera
- El número de entradas en un intervalo de un histograma

# Distribución de Poisson

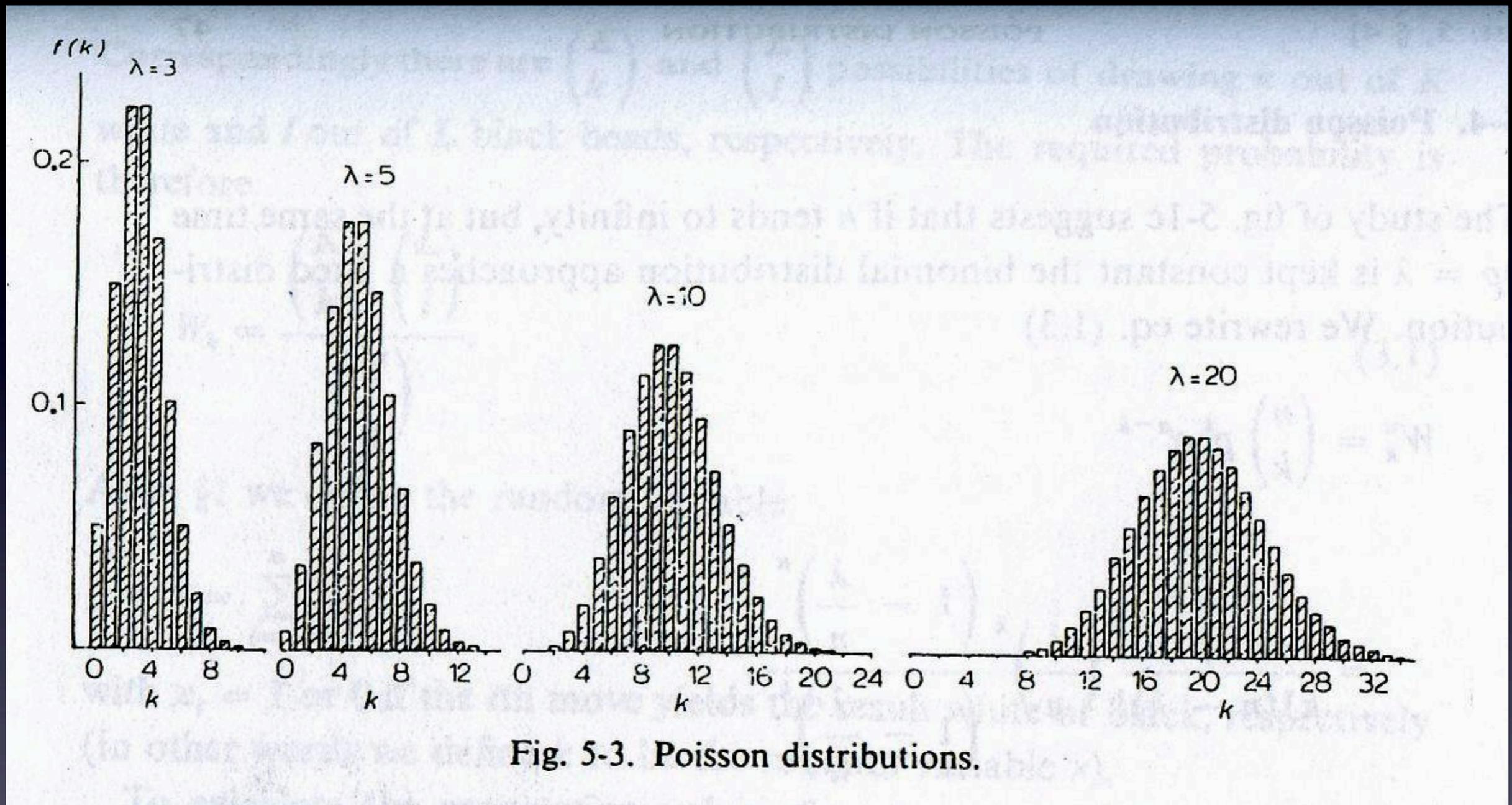


Fig. 5-3. Poisson distributions.

## Comentario:

En general, todos los contajes de sucesos independientes y de baja probabilidad, que se producen a ritmo constante, se supone que se distribuyen según una distribución de Poisson

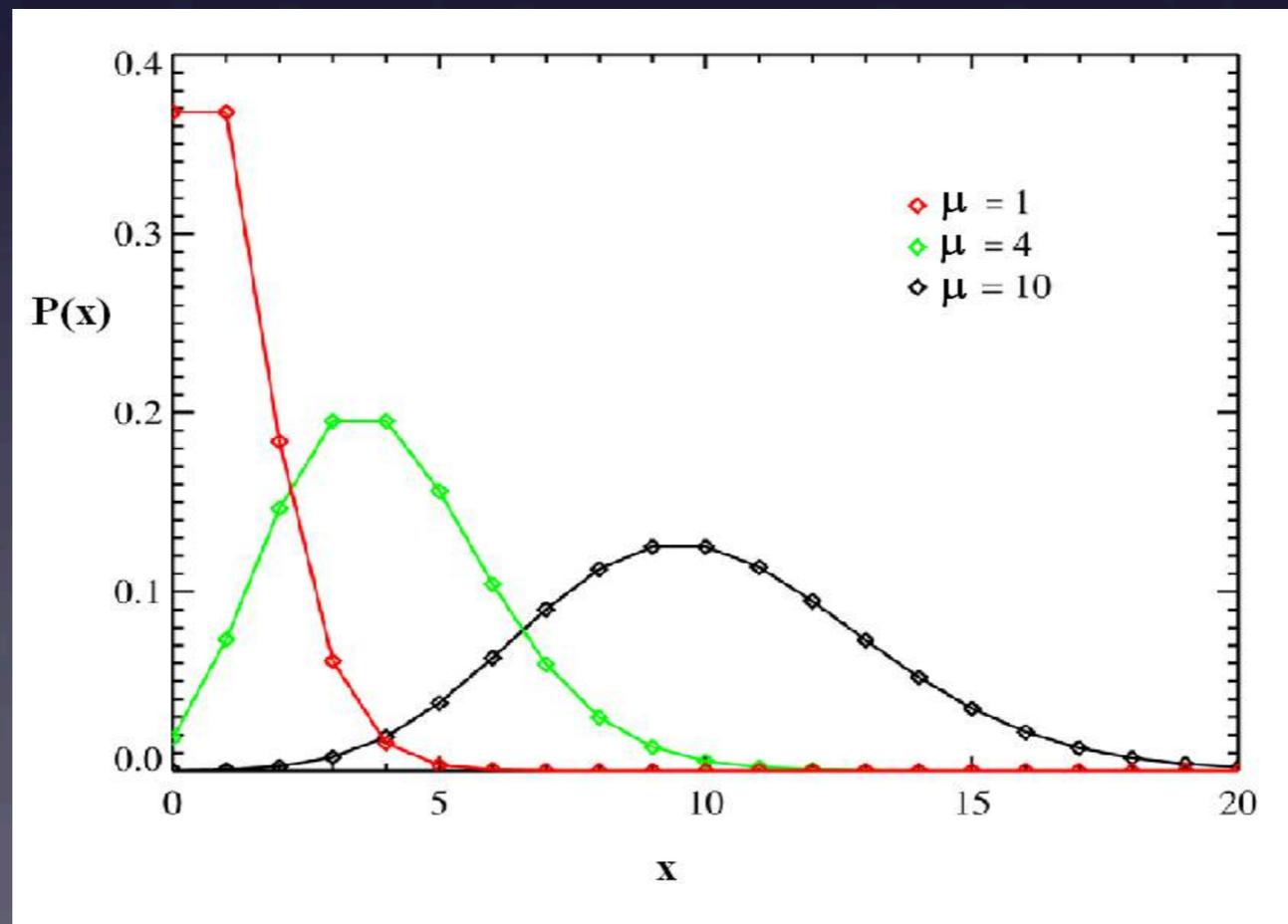
# Distribución de Poisson

La distribución de Poisson se obtiene como límite de la distribución binomial cuando se hace:

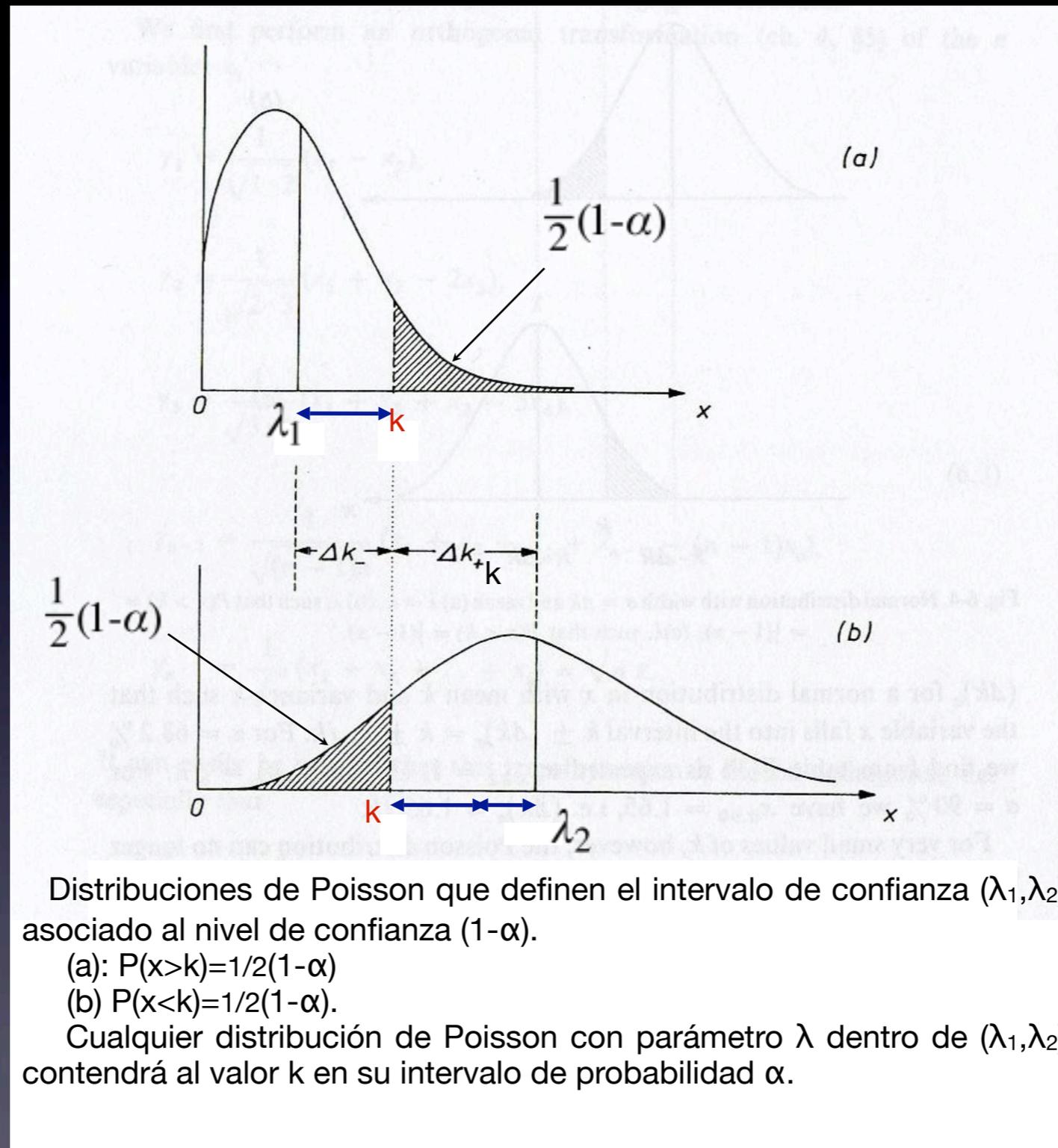
$p \Rightarrow 0$  y  $N \Rightarrow \infty$  manteniendo  $p \cdot N = \text{cte.}$

Requérrese que, en la binomial:  $\mu = N \cdot p$  y  $V = N \cdot p \cdot (1-p)$

Tras el cálculo de los límites:  $\mu = N \cdot p$  y  $V = \mu \cdot (1-p) \Rightarrow V = \mu$  (Poisson)



# Distribución de Poisson: Intervalos de confianza



Ejemplo de intervalo de confianza no simétrico

# Distribución Gaussiana o Normal

Un conjunto de medidas (de una magnitud continua) se agrupan siguiendo una distribución gaussiana si podemos suponer que se dan las siguientes condiciones:

1. Cada una de las medidas está afectada por infinitos factores incontrolables y aleatorios
2. Cada uno de los factores modifica la medida en una cantidad infinitesimal
3. Los factores pueden afectar a la medida por exceso o por defecto con la misma probabilidad

La expresión analítica de una distribución gaussiana depende de **2 parámetros independientes**:  $\mu$  y  $\sigma$ , y toma la siguiente forma:

$$P(x;\mu,\sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

El factor constante que aparece delante de la exponencial es el factor de normalización.

*Nota: La distribución gaussiana no es integrable analíticamente. Curiosamente, sí es integrable la distribución gaussiana bidimensional  $P(x,y)$ .*

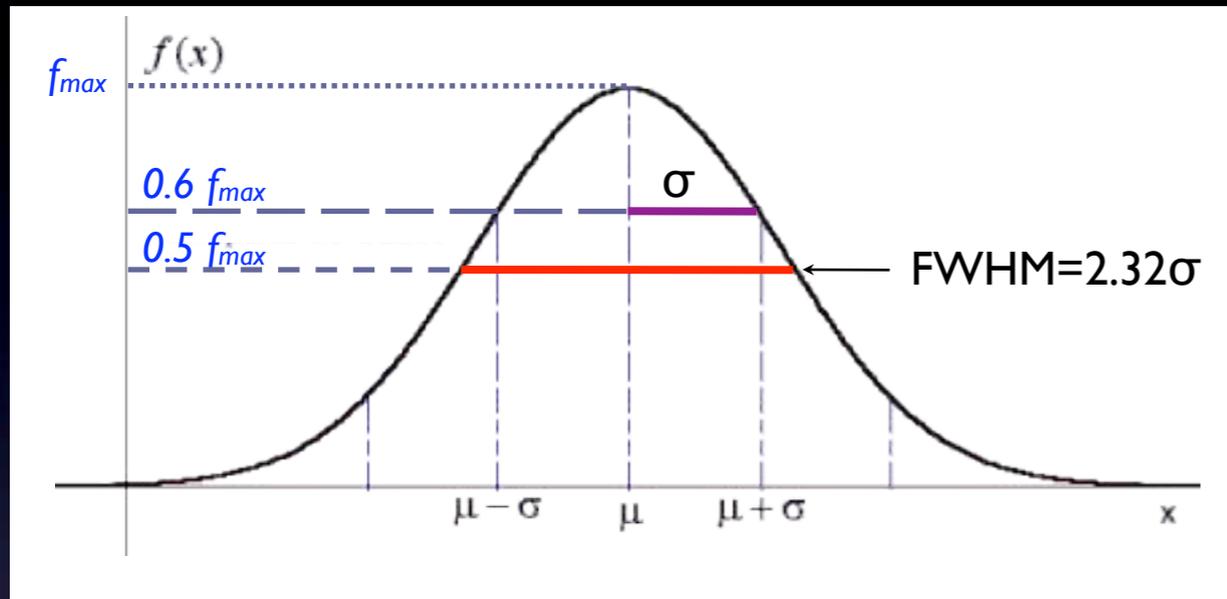
$$\begin{aligned}\mu &= E(k) = \mu \\ V &= \sigma^2\end{aligned}$$

**Comentario:**

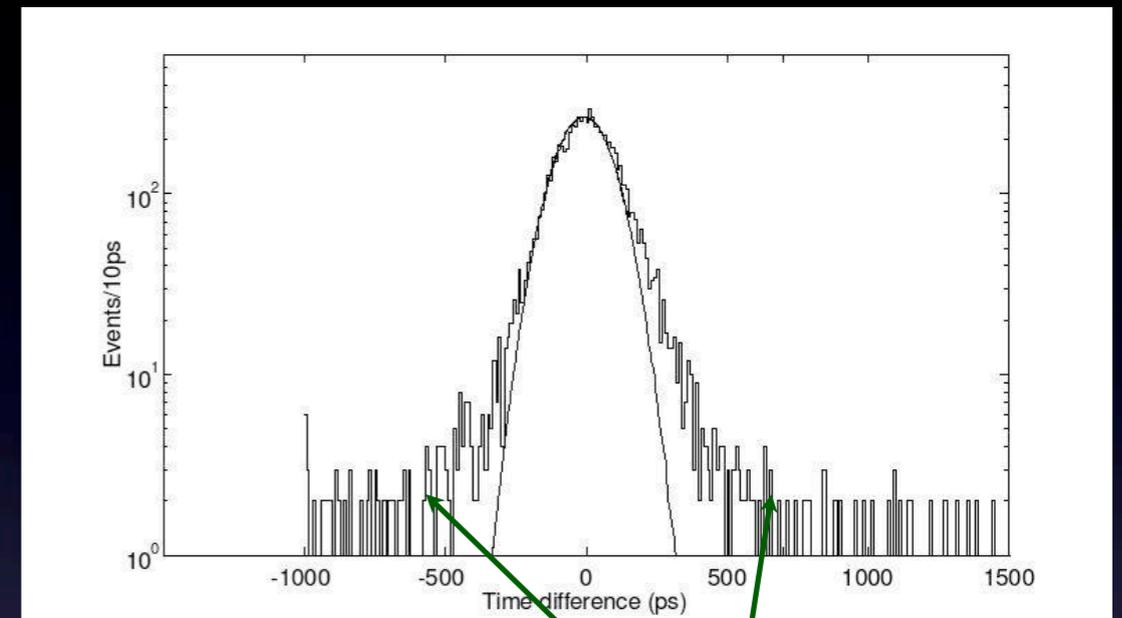
En general, toda las medidas continuas de una magnitud constante, realizadas con suficiente precisión y sin sesgos y sin contribuciones importantes al error, se puede suponer que se agrupan según una distribución gaussiana

# Distribución Gaussiana o Normal

Su representación gráfica es:

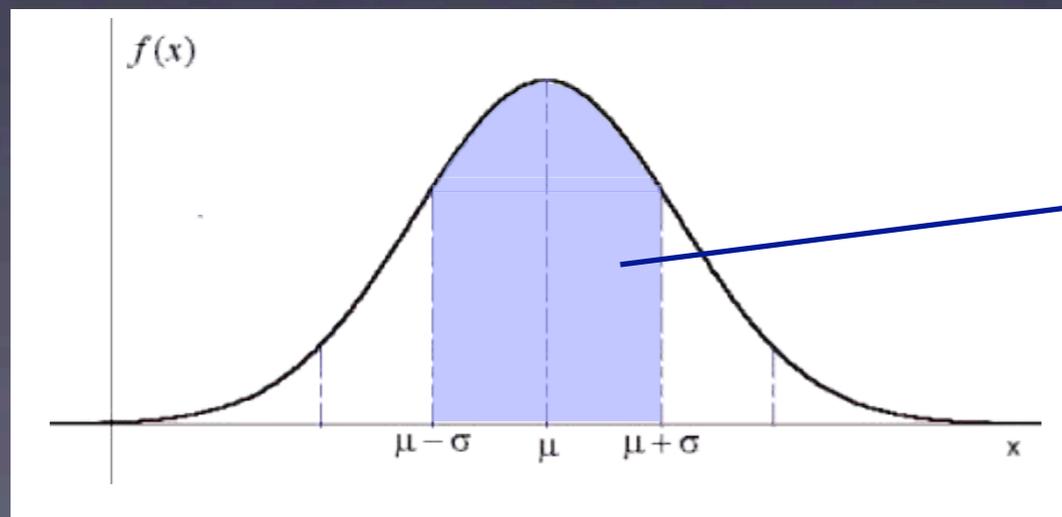


Ejemplo, en escala logarítmica:



Algunas propiedades:

- $f(x=\mu \pm \sigma) = 0.6 \cdot f_{max}$
- La Anchura de la Función a Mitad de la Altura, FWHM (Full Width at Half Maximum) =  $2.32\sigma$
- Valor de la integral entre distintos intervalos de probabilidad:



¡Colas no gaussianas!

Extremo Inferior	Extremo superior	Integral
$-\sigma$	$+\sigma$	68%
$-2\sigma$	$+2\sigma$	95,4%
$-3\sigma$	$+3\sigma$	99,7%
$-4\sigma$	$-4\sigma$	99,99%

# Comparación de distribuciones de Poisson y Gaussiana

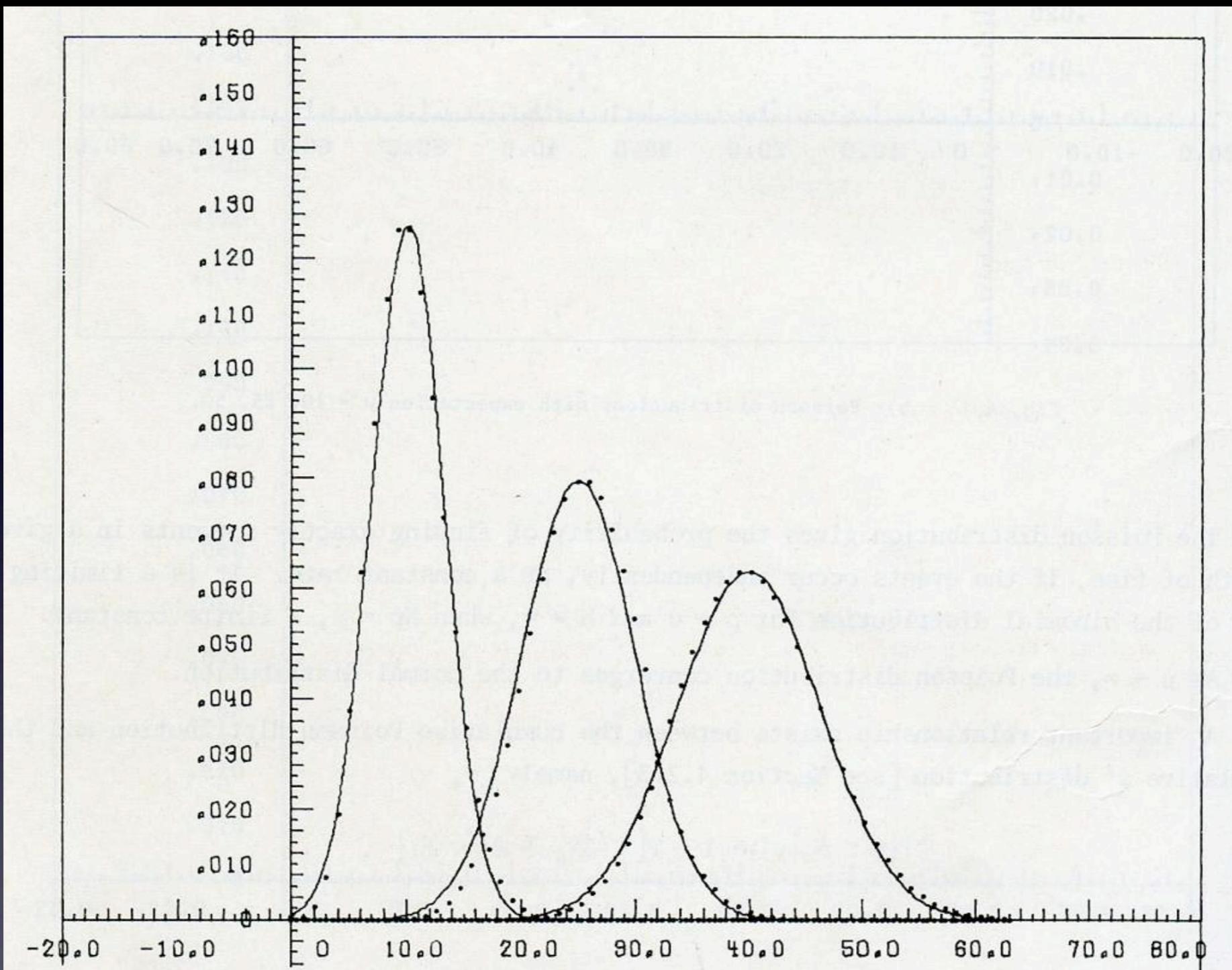


Fig. 4.2 : Comparison of Poisson distributions (dotted) with normal distributions (solid lines) of same mean and variance.

a)

b) Poisson distributions with  $\mu = 10, 25, 40$ . Normal distributions  $N(10,10)$ ,  $N(25,25)$ ,  $N(40,40)$ .

# Distribución Normal Tipificada

La distribución normal centrada de  $\mu=0$  y con anchura  $\sigma=1$  se la llama distribución normal tipificada o distribución normal  $N(0,1)$

$$P(x;0,1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

Los valores de esta distribución son los que habitualmente vienen en las tablas estadísticas o los que producen los generadores de números aleatorios de los ordenadores.

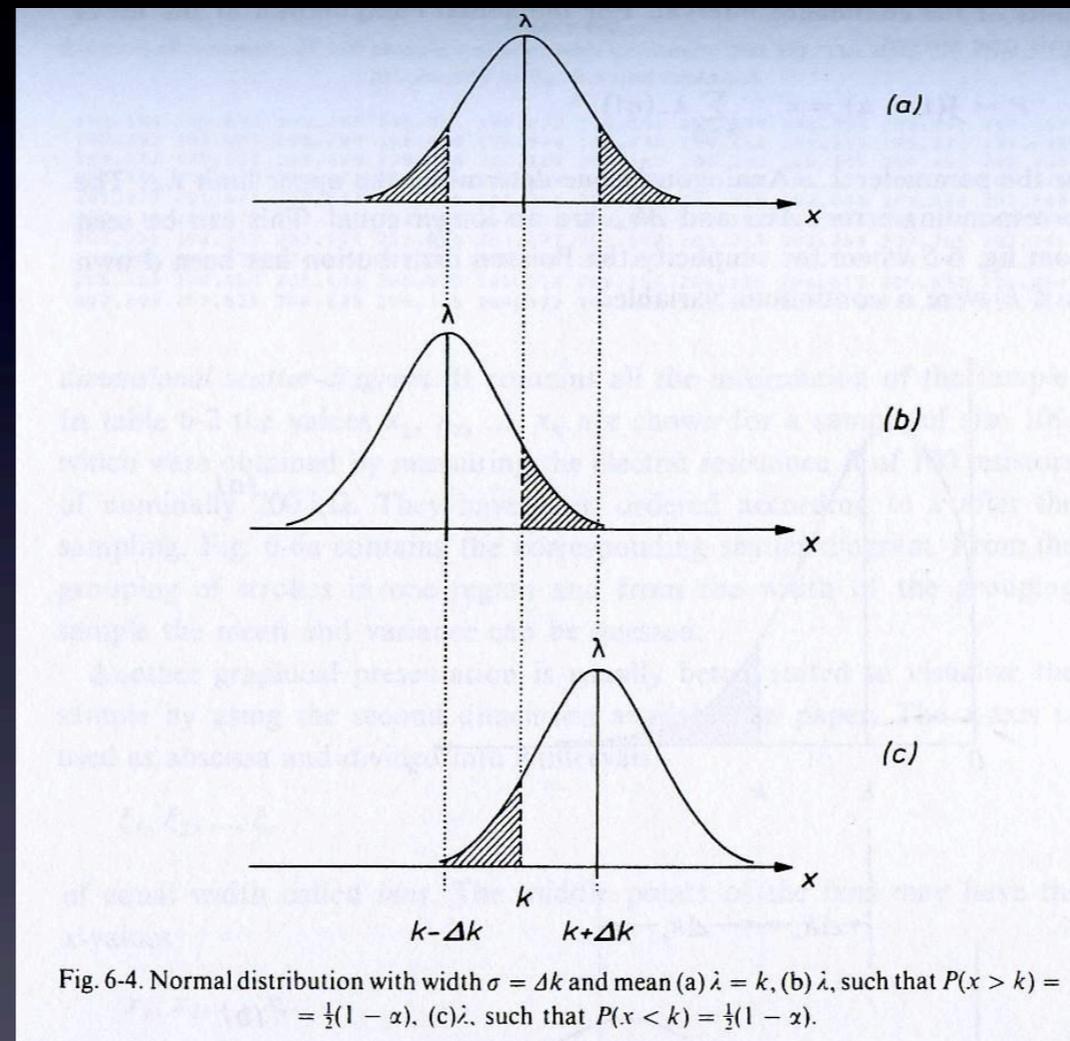
Transformación para pasar de:  $N(0,1) \Rightarrow N(\mu, \sigma)$ :

$$x = \sigma (x_{\text{norm}} + \mu)$$

Transformación para pasar de:  $N(\mu, \sigma) \Rightarrow N(0,1)$ :

$$x_{\text{norm}} = (1/\sigma) (x - \mu)$$

# Distribución gaussiana: Intervalos de confianza



Ejemplo de intervalo de confianza simétrico

# Distribución Chi<sup>2</sup> (o Chi Cuadrado)

Los resultados de la medida de una variable continua se agrupan según un distribución Chi cuadrado, o Chi<sup>2</sup> de **k grados de libertad** si se verifica:

- Dicha variable es la **suma de los cuadrados de k variables con distribución gaussiana N(0,1)**
- Dicha variable es la **suma de los cuadrados de n variables con distribución gaussiana N(0,1)** entre las que existen **n-k ligaduras** o relaciones independientes que las relacionan

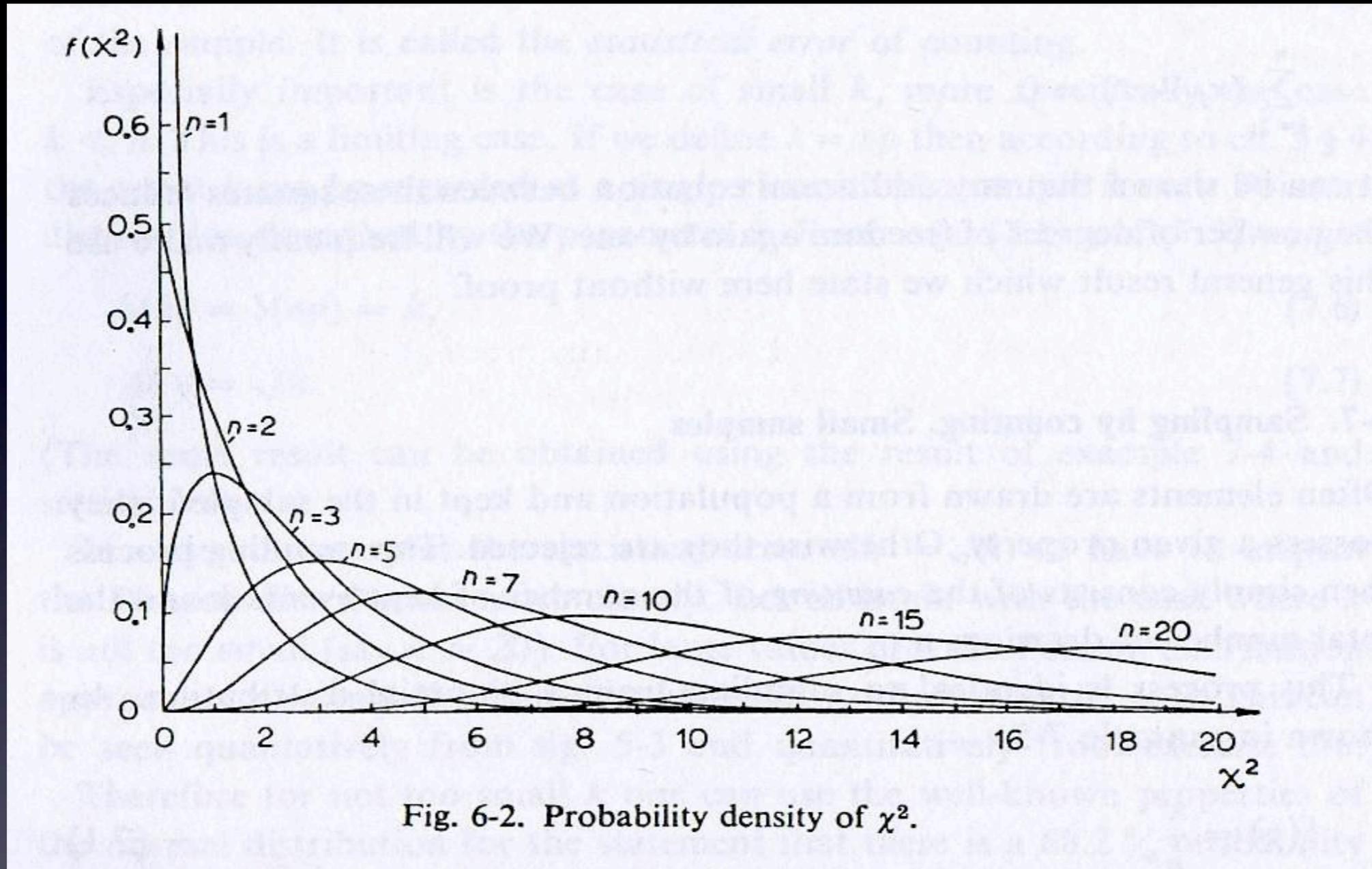
$$f(z;k) = \frac{\left(\frac{1}{2}\right)^{\frac{k}{2}}}{\Gamma\left(\frac{k}{2}\right)} z^{\frac{k}{2}-1} e^{-\frac{z}{2}}$$

- La variable z toma solamente valores positivos. Sus parámetros característicos son:

$$\begin{aligned}\mu &= k \\ V &= 2k\end{aligned}$$

- La distribución de Chi<sup>2</sup> es muy útil en las técnicas de determinación paramétrica

# Chi<sup>2</sup> o Chi Cuadrado



Distribuciones de Chi cuadrado para diferentes grados de libertad, n

# Chi<sup>2</sup> o Chi Cuadrado

Ejemplo:

Sean  $x$ ,  $y$  y  $z$  tres variables que siguen una distribución  $N(0,1)$  y tales que  $x+y+z=cte$ .

Entonces, el resultado de la medida de la variable  $n=x^2+y^2+z^2$  sigue una distribución Chi cuadrado con  $k=3-1=2$  grados de libertad.

Ejemplo:

Sea un conjunto de variables  $x_i$ , que se miden independientemente y que siguen distribuciones normales de la forma  $N(\mu_i, \sigma_i)$ . Entonces, la cantidad,

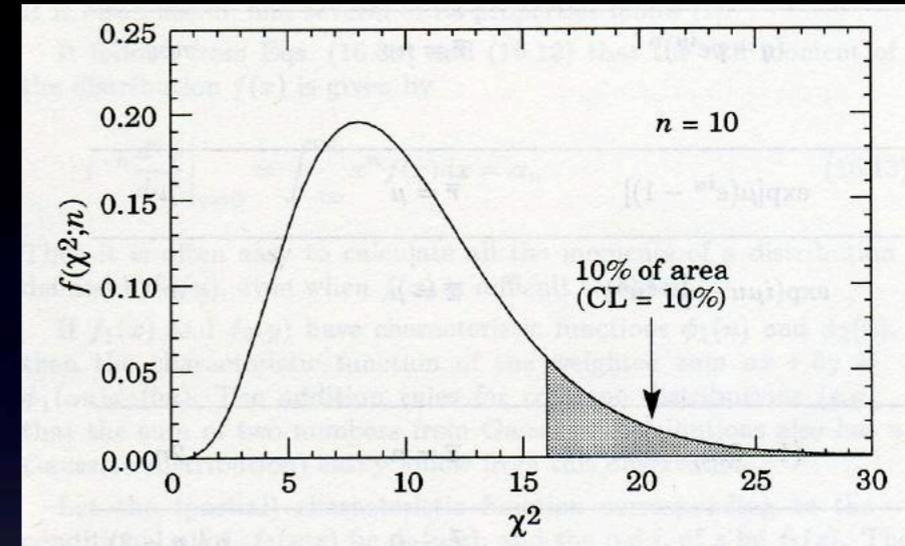
$$s^2 = \sum_{i=1}^n \left( \frac{x_i - \mu_i}{\sigma_i} \right)^2$$

se espera que siga una distribución de Chi<sup>2</sup> con  $n$  grados de libertad.

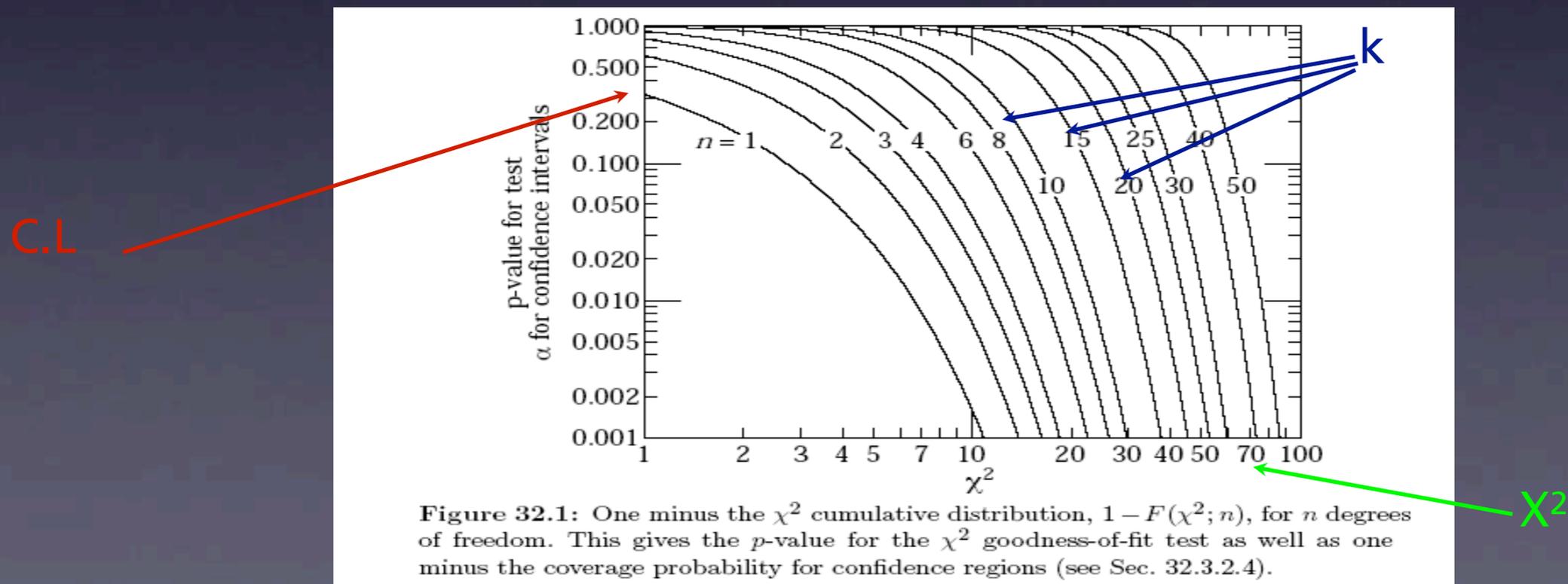
# Chi<sup>2</sup> o Chi Cuadrado

Para la distribución de Chi<sup>2</sup> se define el Nivel de Confianza (o C.L) asociado a un cierto valor de la función X<sup>2</sup>, como:

$$C.L.(\chi^2) = \int_{\chi^2}^{\infty} f(z;k).dz$$



Esta cantidad está tabulada para diferentes valores del grado de libertad  $k$  (y está incorporada a las funciones de la librería de casi todos los programas de análisis estadístico de los ordenadores)



# Resumen de principales distribuciones estadísticas

## 10 31. Probability

**Table 31.1.** Some common probability density functions, with corresponding characteristic functions and means and variances. In the Table,  $\Gamma(k)$  is the gamma function, equal to  $(k-1)!$  when  $k$  is an integer.

Distribution	Probability density function $f$ (variable; parameters)	Characteristic function $\phi(u)$	Mean	Variance $\sigma^2$
Uniform	$f(x; a, b) = \begin{cases} 1/(b-a) & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$	$\frac{e^{ibu} - e^{iau}}{(b-a)iu}$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
Binomial	$f(r; N, p) = \frac{N!}{r!(N-r)!} p^r q^{N-r}$ $r = 0, 1, 2, \dots, N; \quad 0 \leq p \leq 1; \quad q = 1 - p$	$(q + pe^{iu})^N$	$Np$	$Npq$
Poisson	$f(n; \nu) = \frac{\nu^n e^{-\nu}}{n!}; \quad n = 0, 1, 2, \dots; \quad \nu > 0$	$\exp[\nu(e^{iu} - 1)]$	$\nu$	$\nu$
Normal (Gaussian)	$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp(-(x-\mu)^2/2\sigma^2)$ $-\infty < x < \infty; \quad -\infty < \mu < \infty; \quad \sigma > 0$	$\exp(i\mu u - \frac{1}{2}\sigma^2 u^2)$	$\mu$	$\sigma^2$
Multivariate Gaussian	$f(\mathbf{x}; \boldsymbol{\mu}, V) = \frac{1}{(2\pi)^{n/2} \sqrt{ V }}$ $\times \exp[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T V^{-1}(\mathbf{x} - \boldsymbol{\mu})]$ $-\infty < x_j < \infty; \quad -\infty < \mu_j < \infty; \quad  V  > 0$	$\exp[i\boldsymbol{\mu} \cdot \mathbf{u} - \frac{1}{2}\mathbf{u}^T V \mathbf{u}]$	$\boldsymbol{\mu}$	$V_{jk}$
$\chi^2$	$f(z; n) = \frac{z^{n/2-1} e^{-z/2}}{2^{n/2} \Gamma(n/2)}; \quad z \geq 0$	$(1 - 2iu)^{-n/2}$	$n$	$2n$
Student's $t$	$f(t; n) = \frac{1}{\sqrt{n\pi}} \frac{\Gamma[(n+1)/2]}{\Gamma(n/2)} \left(1 + \frac{t^2}{n}\right)^{-(n+1)/2}$ $-\infty < t < \infty; \quad n$ not required to be integer	—	0 for $n \geq 2$	$n/(n-2)$ for $n \geq 3$
Gamma	$f(x; \lambda, k) = \frac{x^{k-1} \lambda^k e^{-\lambda x}}{\Gamma(k)}; \quad 0 < x < \infty;$ $k$ not required to be integer	$(1 - iu/\lambda)^{-k}$	$k/\lambda$	$k/\lambda^2$

# Test de hipótesis: Determinación paramétrica

- Cuando unas medidas siguen una cierta distribución matemática conocida, sus parámetros característicos se pueden determinar a partir de las medidas características (valor medio, varianza...) de la distribución de los datos
- Cuando las medidas siguen una distribución desconocida, para la que se puede tener modelo o no, el problema es más complicado y hay que recurrir a las técnicas de **test de hipótesis** o de **determinación paramétrica**:

Consiste en determinar el valor de aquellos **parámetros** que caracterizan una cierta función que se supone que interpretan los datos

P. ej: Si se sospecha que unos datos  $\{y_i\}$  siguen un comportamiento lineal con una variable  $x$ , que hemos variado durante la medida, de la forma:

$$y_i = a \cdot x_i + b$$

se pretende determinar cuáles son los valores de los **parámetros  $a$  y  $b$**  que mejor interpretan los datos observados y asignarles, además, una incertidumbre

Los principales métodos de determinación paramétrica son:

- Método de mínimo Chi cuadrado
- Método de máxima verosimilitud

# Método del mínimo $\chi^2$ para la determinación paramétrica

Sea el caso en que se dispone de:

- un conjunto de  $n$  medidas  $\{y_i\}$  de la variable  $y$ , con varianzas respectivas  $\sigma_i^2$  realizadas respectivamente para diferentes valores  $\{x_i\}$  de la variable  $x$  y
- se dispone de modelo, o teoría, por la cual se debe de verificar que:

$$y = f(x; \mathbf{a}) \quad \text{siendo } \mathbf{a} = \{k \text{ parámetros desconocidos}\}$$

Si, efectivamente, los datos obedecen a dicha relación:

- cada  $y_i$  es una medida del valor  $f(x_i; \mathbf{a})$  afectada por la incertidumbre  $\sigma_i$  y se espera, entonces, que

- las cantidades: 
$$\frac{y_i - f(x_i; \mathbf{a})}{\sigma_i}$$

deben seguir, cada una de ellas, una distribución  $N(0,1)$ .

Entonces, la cantidad  $s = \sum \left( \frac{y_i - f(x_i; \mathbf{a})}{\sigma_i} \right)^2$  debe de seguir una distribución  $\chi^2$

# Método del mínimo Chi<sup>2</sup> para la determinación paramétrica

A partir de la función:

$$s = \sum \left\{ \frac{y_i - f(x_i; \mathbf{a})}{\sigma_i} \right\}^2$$

el método de mínimo Chi<sup>2</sup> consiste en calcular aquel conjunto de valores de los **k parámetros desconocidos** que hacen mínimo el valor de **s**.

Seguidamente se puede buscar el valor de **s** en la tabla de Chi<sup>2</sup> y determinar el Nivel de Confianza asociado (que es la probabilidad de que hubieramos obtenido la observación experimental **{y<sub>i</sub>}** si la relación propuesta  $y = f(x_i; \mathbf{a})$  fuese cierta)

## Procedimiento:

El cálculo de la familia de parámetros **a** que minizan el valor de **s** se puede realizar por:

- Mediante un procedimiento analítico (este método es muy sencillo para polinimios de cualquier orden)
- Mediante un método de minimización numérico con el ordenador (imprescindible para funciones complicadas)

# Método del mínimo Chi<sup>2</sup>

Interpretación de los resultados:

Sean unas medidas sobre las que se han hecho varios ajustes polinómicos por el método de mínimos cuadrados. Se desea saber cuál la mejor estimación de parámetros.

Se disponen de 8 medidas y se han probado polinomios de orden 2 ( $y=a+bx$ ), 3 ( $y=a+bx+cx^2$ ) y 4 ( $y=a+bx+cx^2+dx^3$ ). Los resultados de los ajustes se muestran en la tabla:

Número de <sup>N</sup> puntos	Número de <sup>n:</sup> parámetros	núm. de grados de libertad <sup>k:</sup>	Chi <sup>2</sup>	Chi <sup>2</sup> /k	CL
8	2	6	12.3	2.1	
8	3	5	7.2	1.4	
8	4	4	6.4	1.6	

La mejor elección es la de un polinomio cuadrático,  $n=3$ .

Al añadir más parámetros, la calidad del ajuste mejora pero se pierde significación estadística

# Método del mínimo $\chi^2$

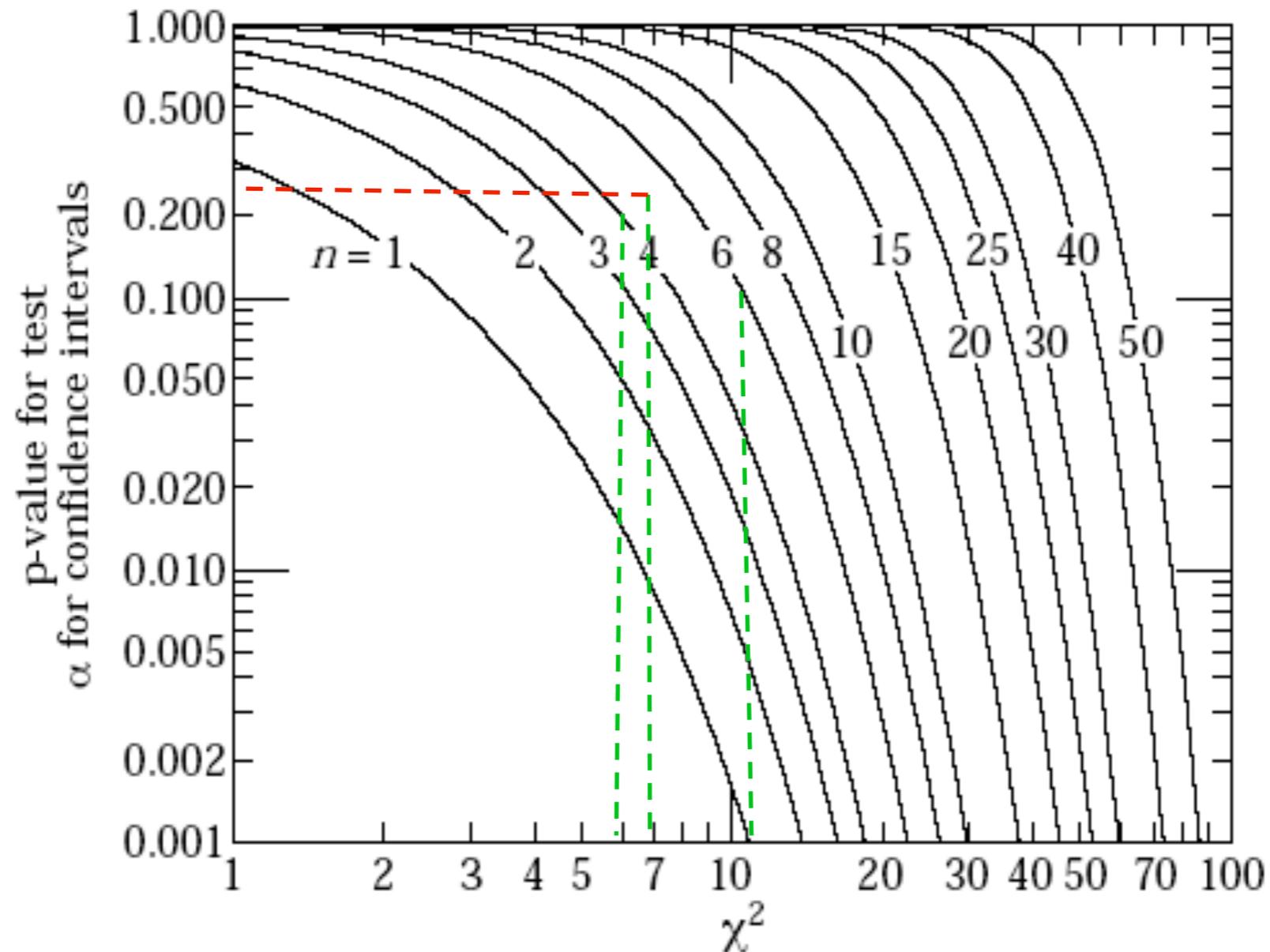


Figure 32.1: One minus the  $\chi^2$  cumulative distribution,  $1 - F(\chi^2; n)$ , for  $n$  degrees of freedom. This gives the  $p$ -value for the  $\chi^2$  goodness-of-fit test as well as one minus the coverage probability for confidence regions (see Sec. 32.3.2.4).

# Una historia real

En 1974, un equipo de los laboratorios SLAC y LBL estaba realizando un estudio rutinario de la sección eficaz total electrón-positrón en función de la energía en el centro de masas, para energías de unos pocos GeVs.

La sección eficaz era mas o menos constante, con valores de alrededor de unos  $23 \pm 3$  nb. Alrededor del valor de unos 3.2GeV la sección eficaz tomaba un valor de unos 30nb. Este valor estaba a unas 2 desviaciones típicas por encima del valor medio que se estaba observando, comportamiento relativamente esperable cuando se miden muchos puntos (la probabilidad de que una desviación así tenga lugar es del orden de un 5%). A un investigador poco cuidadoso no le habría llamado la atención y habría pasado de largo sobre dicho punto.

Sin embargo, a aquellos investigadores sí les llamó la atención por lo que decidieron aumentar la escala de la representación y analizar con mayor detalle dicho efecto. Como consecuencia descubrieron un nuevo hadrón que no encajaba en el modelo de tres quarks existente entonces. Denominaron a dicha partícula  $\psi$  que, con una masa de  $3.105 \pm 0.003$ GeV, resultó ser un estado ligado de un quark encantado con su correspondiente antiquark. (La misma partícula se descubrió casi simultáneamente en el laboratorio BNL, al otro lado de los EEUU, donde la llamaron partícula J).

En 1976 el premio Nobel de Física fue concedido de forma compartida a los investigadores que dirigían los dos equipos responsables del descubrimiento: Burton Richter (SLAC) y Samuel Ting (BNL)

Moraleja:

El premio Nobel puede estar merodeando alrededor de lo que parece un dato sin interés. Por lo tanto, investiga (tomando mas medidas) cualquier dato extraño antes de despreciarlo.

Ref. Daryl W. Preston: Experiments in Physics (J. Wiley & Sons)

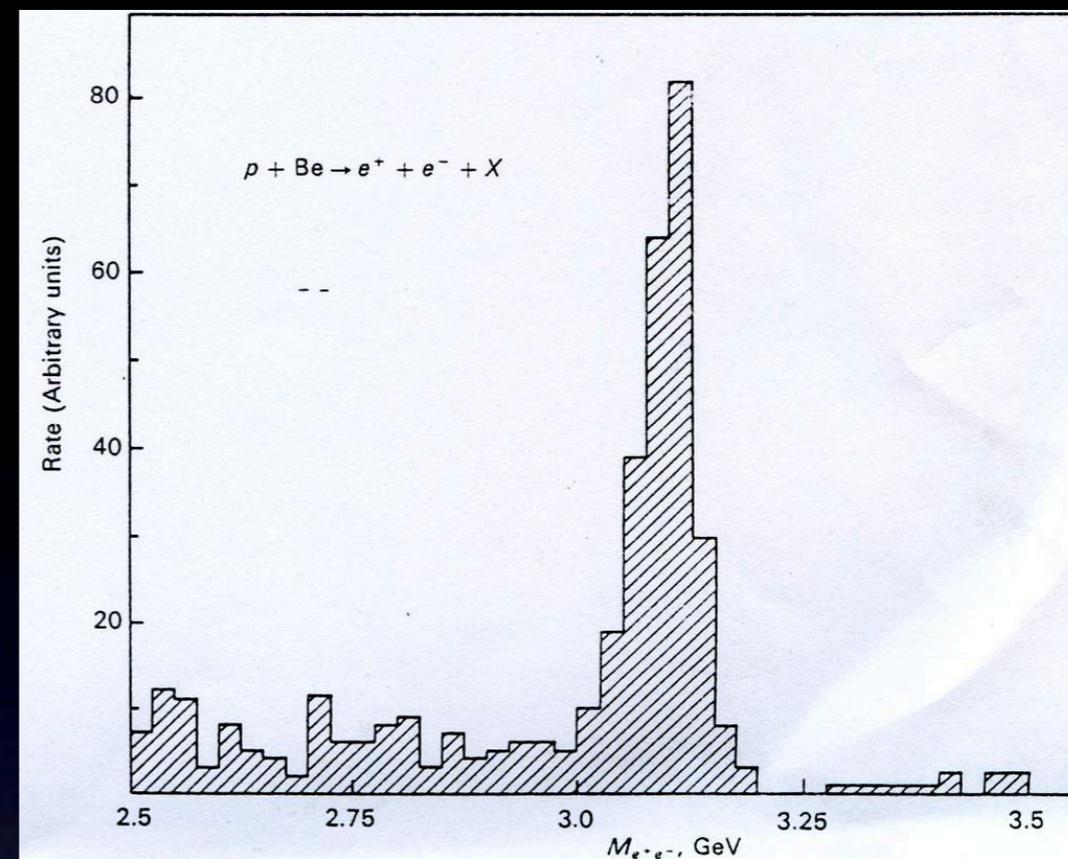


Fig. 5.11 Results of Aubert *et al.* (1974) indicating the narrow resonance  $\psi/J$  in the invariant mass distribution of  $e^+e^-$  pairs produced in inclusive reactions of protons with a beryllium target. The experiment was carried out with the 28-GeV AGS at Brookhaven National Laboratory.

