

A photograph of two lion cubs in a savanna setting. One cub is on the left, facing right, and the other is on the right, facing left. They are both on their hind legs, reaching towards each other with their front paws. The background is a blurred green field with some tall grass. The text 'AST0212 – 2016-1' is overlaid in large yellow letters across the middle of the image.

AST0212 – 2016-1

Introducción al análisis de datos

Instituto de Astrofísica

Facultad de Física

Pontificia Universidad Católica de Chile



Equipo docente:

Profesor: Alejandro Clocchiatti

Ayudantes:

Francisco Aros (TM6)

Nicolás Castro (TL4)

TM6: Tutoría del martes en módulo 6

TL4: Tutoría del lunes en módulo 4

Nuestro Semestre 2016-1

AST0212				C0 ✓		
Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
6 Mar 2016 Semana 1					C1 ✓	
13 Semana 2	TL1	TM1			C2 ✓	← Control 1 Reparto Tarea 1
20 Semana 3	TL2	TM2			Feriado	
27 Semana 4	TL3	TM3			C3 ✓	
3 Semana 5	TL4	TM4			C4 ✓	
10 Semana 6	TL5	TM5			C5	← Control 2
17 Semana 7	TL6	TM6			C6 – SM1	← Reparto T2
24 Semana 8	TL7	← Entrega Tarea 1			C7 – SM2	
1 May Semana 9	TL8	TM8			C8 – SM3	
8 Semana 10	TL9	← Entrega Tarea 2			C9 – SM4	
15 Semana 11	TL10	TM10			C10	
22 Semana 12	TL11	TM11			C11	
29 Semana 13	TL12	TM12	1 Jun		Feriado	
5 Semana 14	TL13	TM13			C12	
12 Semana 15	TL14	TM14			C13	
19 Tutorías día lunes Módulo 4: Polás Castro						
		Tutorías día martes Módulo 6: Francisco Aros				
					Notas	

Clase previa (Clase 4):

REPASO

1. Temas pendientes

- 1. Observaciones desde Santa Martina

- 1. Herramienta Linux de selección de datos en archivos organizados en columnas: *awk*

2. Breve repaso de la clase previa

- 1. Visualización cualitativa de histogramas.

- 2. Histogramas y funciones de distribución de probabilidad.

- 3. Uso de la FDP para calcular parámetros de la distribución.

Temas del día: Segunda vuelta sobre FDP constante. Otras FDP que hay que conocer: Poisson y Gauss. Modelos de la realidad, distribución subyacente. Test modelo vs. realidad.

Esta clase (Clase 5):

1. Herramienta Linux de selección de datos en archivos organizados en columnas: *awk*
2. Repaso de temas críticos de la clase previa
 1. FDPs que hay que conocer: Constante, Poisson y Gauss.
 2. Modelos de la realidad, distribución subyacente.
 3. Test modelo vs. realidad: χ^2 explicado.
 4. FDP de χ^2 .
3. Viaje sin escalas a la propagación de errores.
4. Correlación.

FDP e histogramas de histogramas REPASO

10^6 números. FDP cte. entre 0 y 1.

$\bar{x} = 0.5001$, $\sigma_x = 0.2886$, bin=0.0001

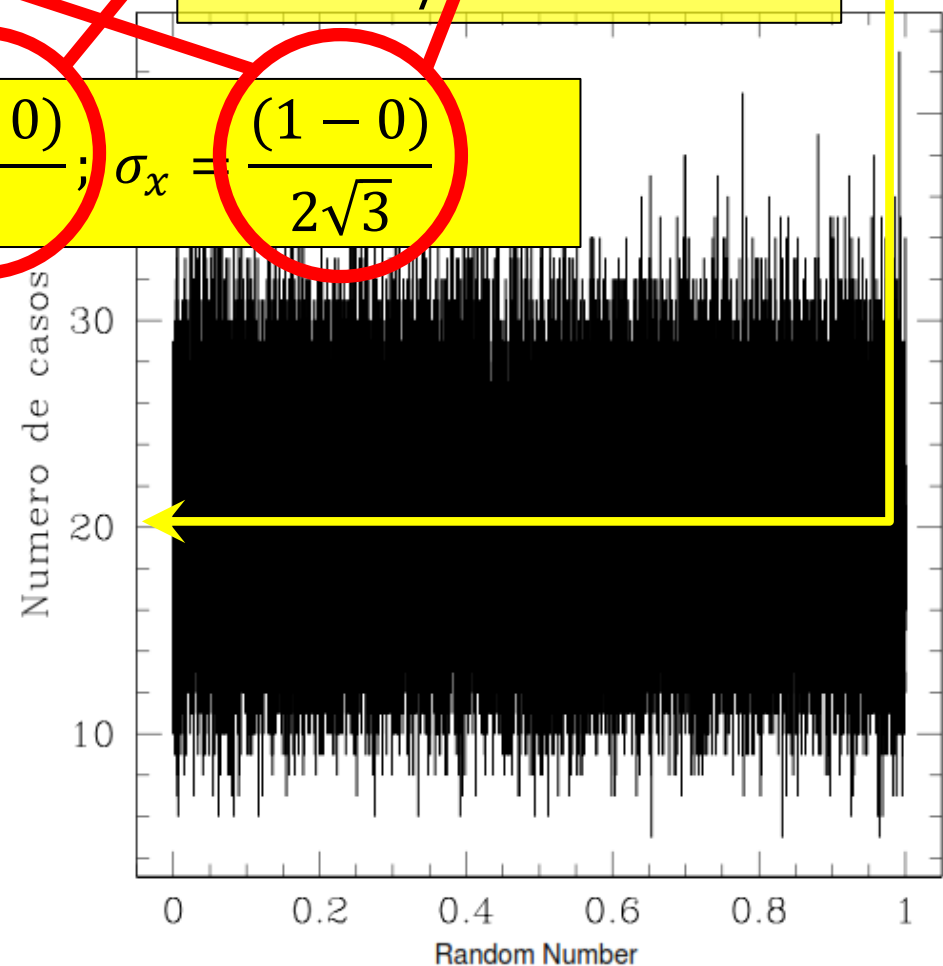
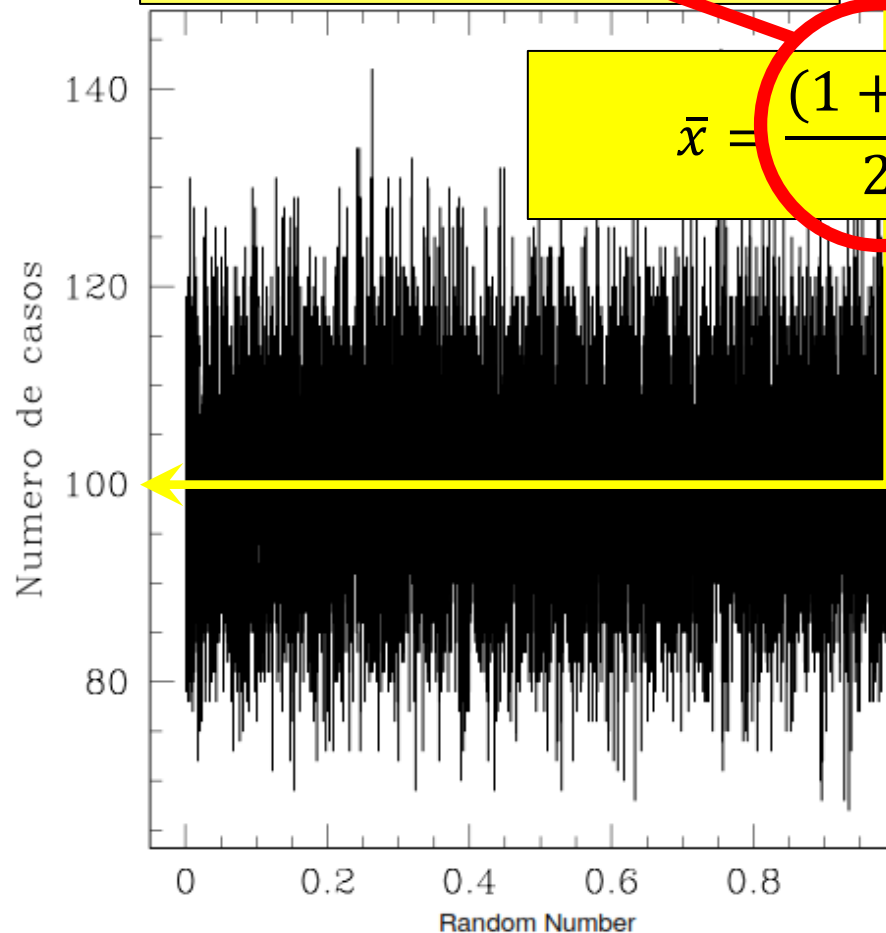
10^6 números. FDP cte. entre 0 y 1.

$\bar{x} = 0.5001$, $\sigma_x = 0.2886$, bin=0.00002

$$1 \times 10^6 / 1 \times 10^4 = 100$$

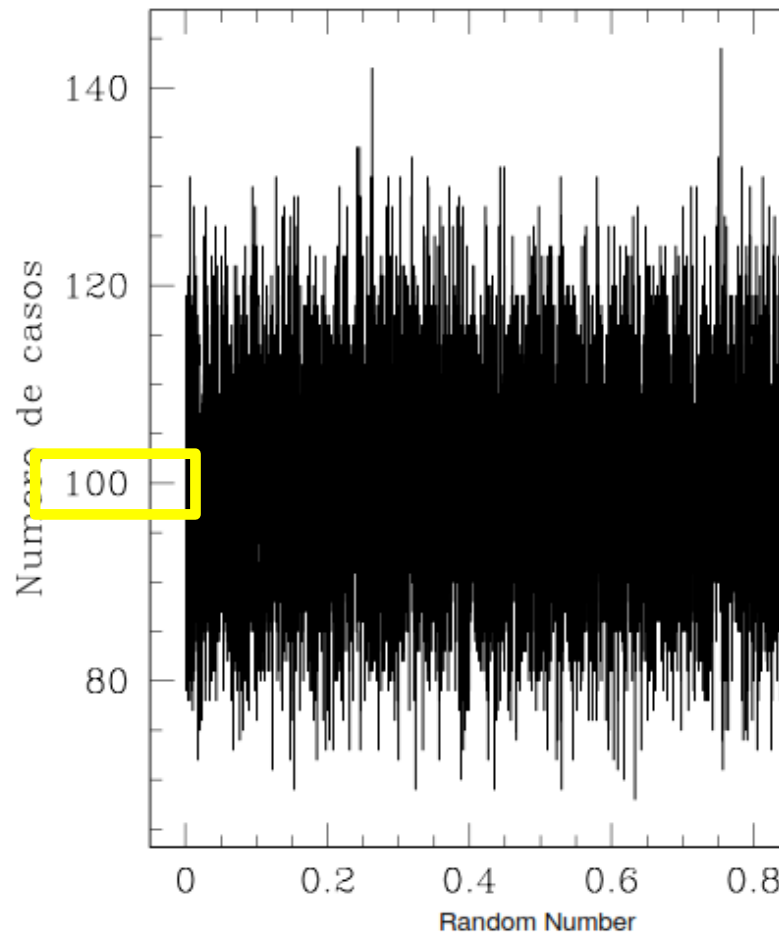
$$1 \times 10^6 / 5 \times 10^4 = 20$$

$$\bar{x} = \frac{(1 + 0)}{2}; \sigma_x = \frac{(1 - 0)}{2\sqrt{3}}$$

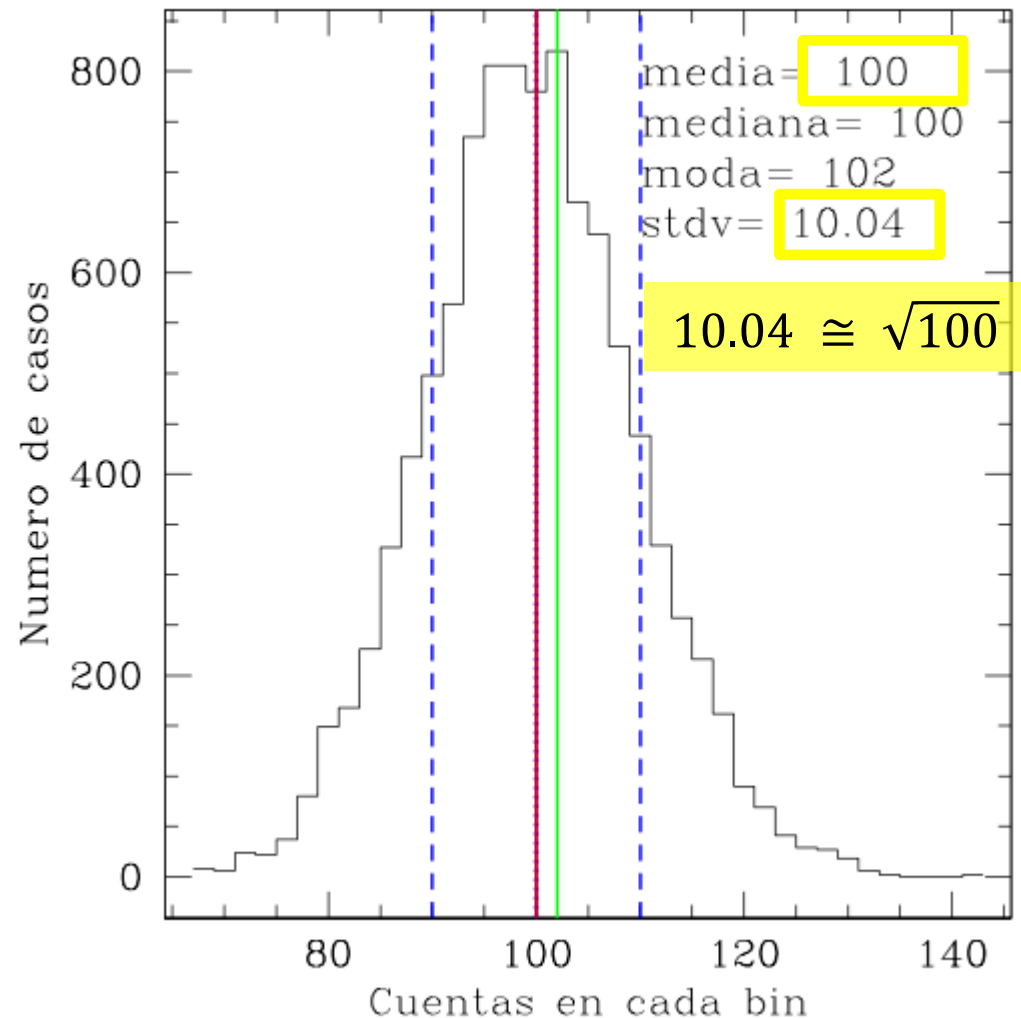


FDP e histogramas de histogramas REPASO

1e6 RdN - Bin 1e-4



hist_values_rdn_1e6.dat2 ; $N_T = 9999$; Bin = 2

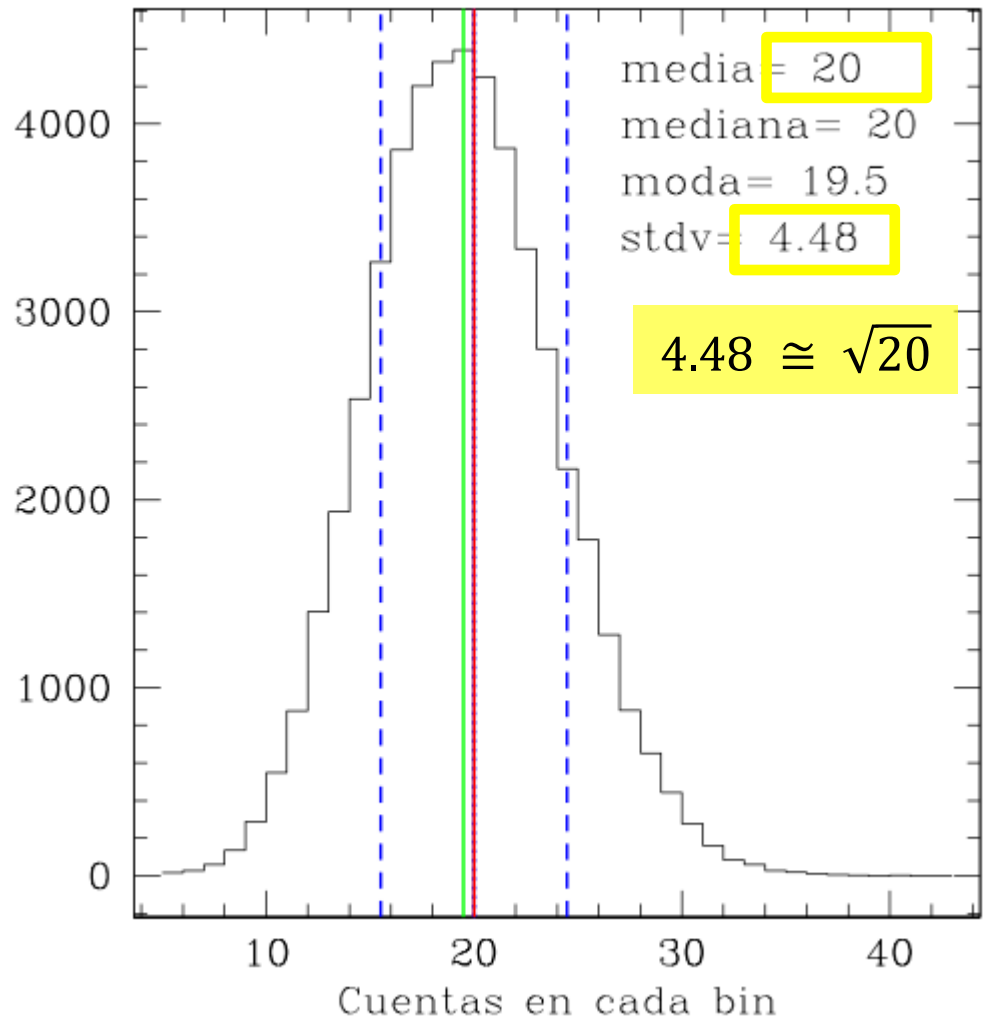
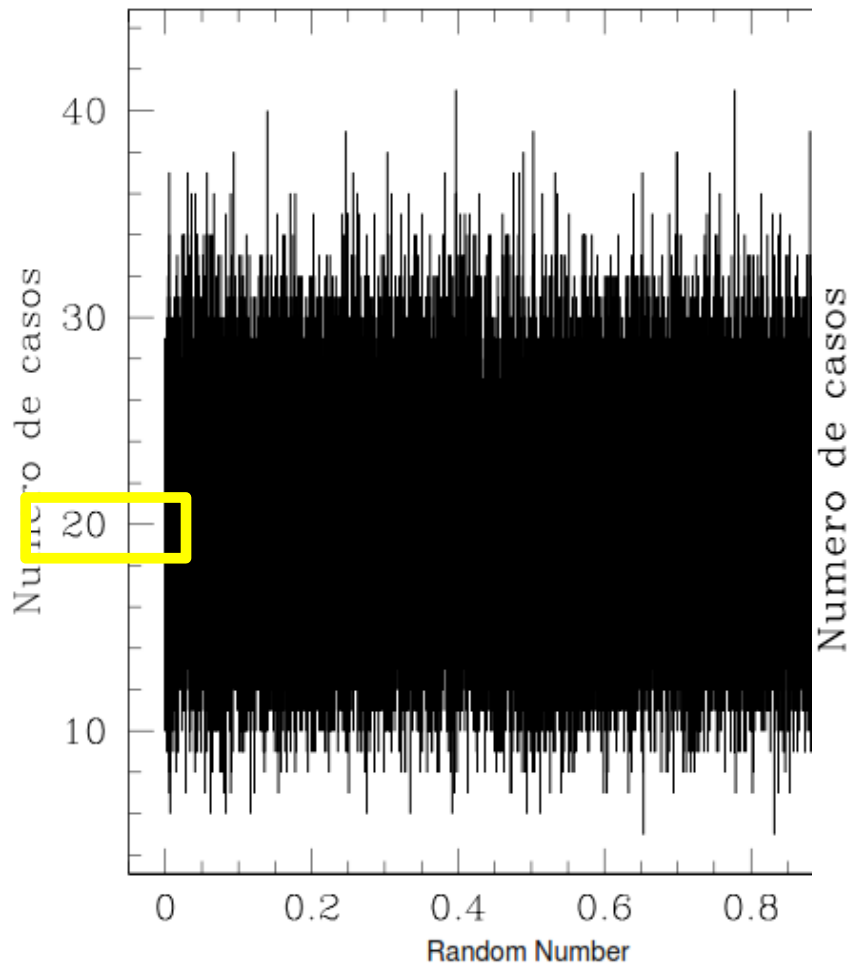


FDP e histogramas de histogramas REPASO

Si clasifico a los números aleatorios en bins más chicos, la FDP que obtengo será la misma, con parámetros diferentes:

1e6 RdN - Bin 2e-5

hist_values_rdn_1e6_2e-5.dat2 ; $N_T = 49999$; Bin = 1



FDP de Poisson

REPASO

La FDP que está detrás de todo esto es la llamada *Distribución de Poisson*, que resulta de contar eventos que suceden en un intervalo (de tiempo o espacio) dado, definido, cuando la probabilidad individual de cada evento es muy baja. Por ejemplo:

1. Decaimiento radioactivo de núcleos atómicos por segundo.
2. Explosiones de SN en un volumen del universo en un intervalo de tiempo.
3. Cantidad de gotas de lluvia que caen en un vaso en un intervalo de tiempo.
4. Número de fotones que llegan a un pixel de un CCD en una exposición.
5. Cantidad de números aleatorios que caen en un bin específico.

La FDP de Poisson, está dada por:

$$P_{\mu}(\nu) = e^{-\mu} \frac{\mu^{\nu}}{\nu!} ; \text{ con } \mu > 0$$

que es, específicamente, la probabilidad de contar ν eventos en el intervalo dado (la ecuación anterior está normalizada).

Puede mostrarse que para esta FDP $\bar{\nu} = \mu$ y $\sigma_{\nu}^2 = \mu$, o sea $\sigma_{\nu} = \sqrt{\mu}$.

Entonces, si la tasa de ocurrencia es R (probabilidad del evento por unidad de intervalo), entonces $\mu = RT$, donde T es el largo del intervalo. Estas ecuaciones aclaran todas las coincidencias anteriores.

FDP de Poisson

REPASO

Esta clase de análisis provee una herramienta muy buena para testear el software que estamos usando y asegurarnos que hace lo que dice que hace:

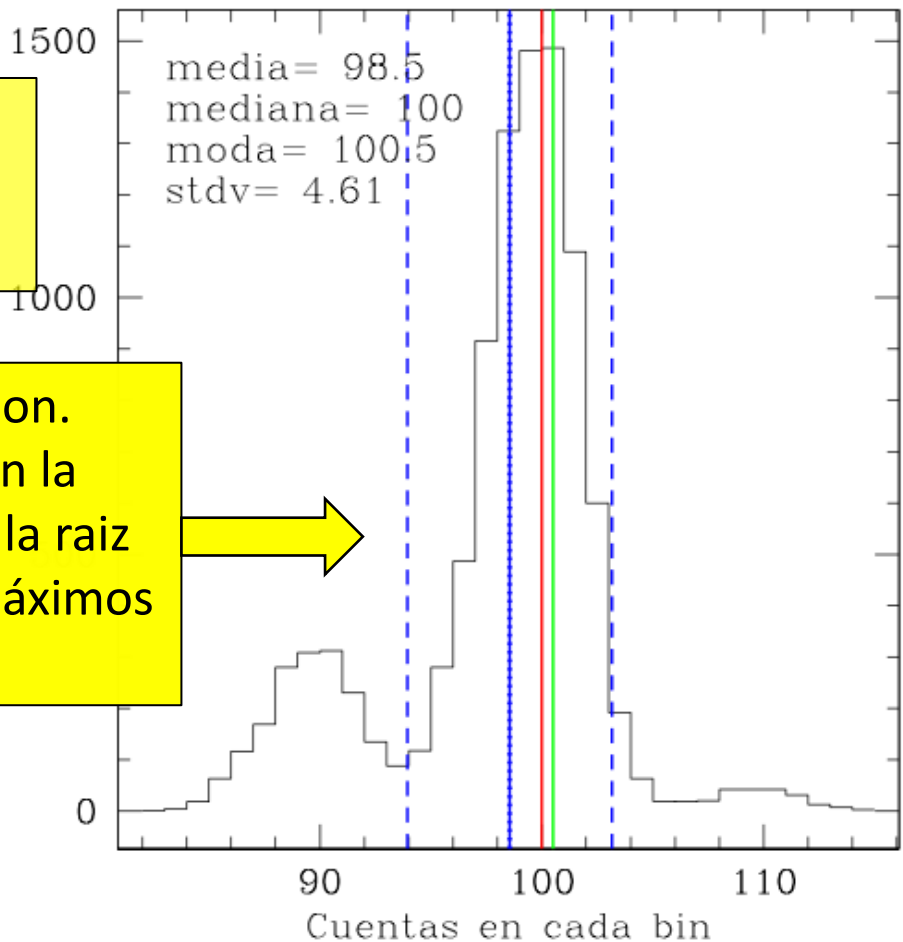
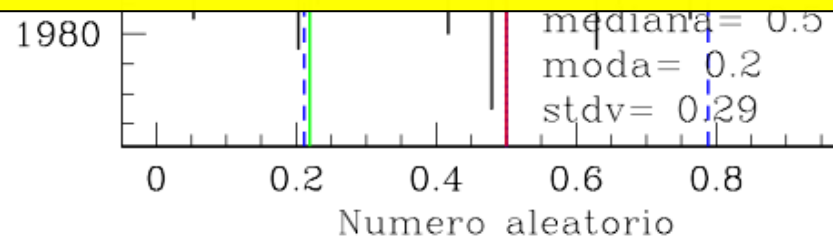
¡No todos los generadores de números al azar que andan por ahí son buenos!

rdn₁e6.dat ; $N_T = 1000000$; Bin = 0.002

hist_values_rdn₁e6.dat ; $N_T = 10000$; Bin = 1

Generador de números al azar RAN1 de
Numerical Recipes (Press, Flannery,
Teukolsky & Vetterling, 1989)

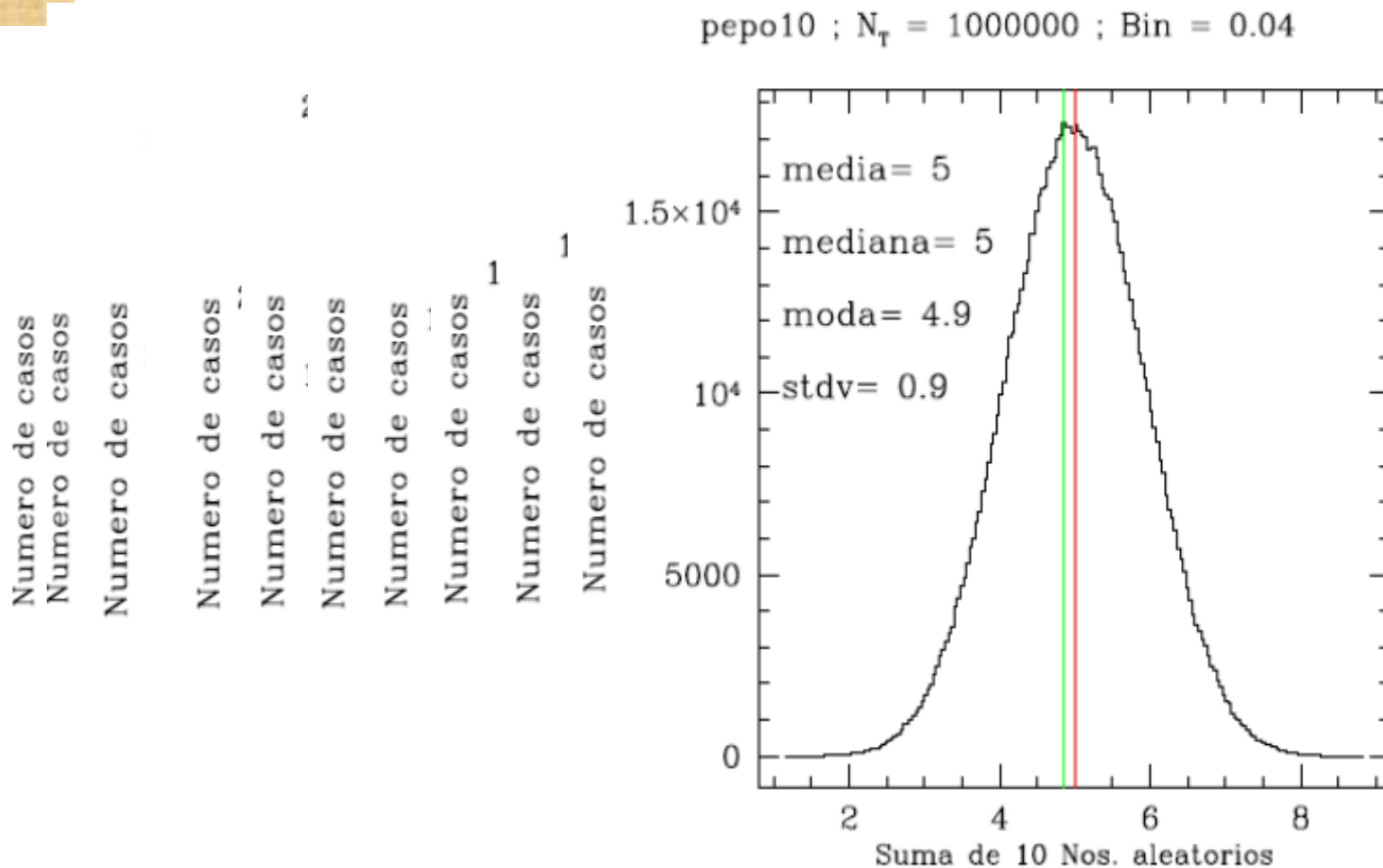
Esta distribución no se parece a la de Poisson.
Tiene un máximo centrado más o menos en la
posición correcta, pero la dispersión no es la raíz
cuadrada del valor medio y muestra dos máximos
secundarios centrados cerca de 90 y 110.



FDP de Gauss

REPASO

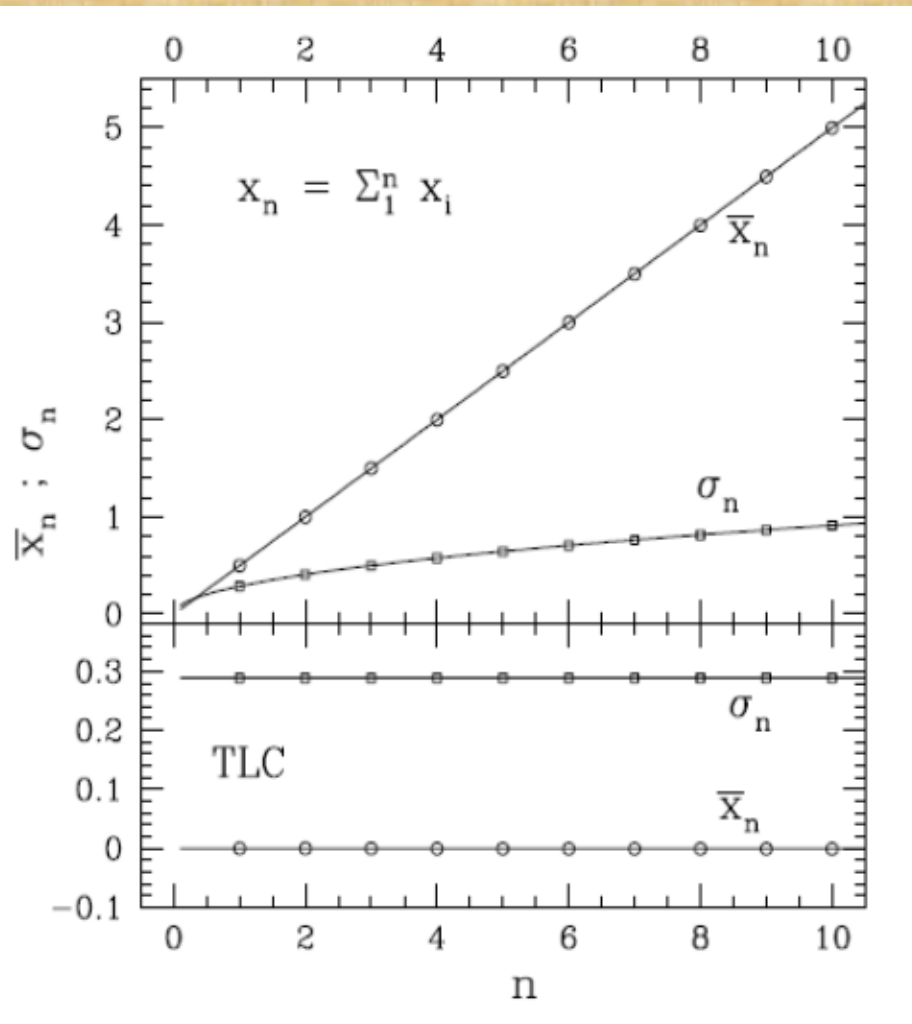
Una forma de presentar la FDP de Gauss podría ser “la distribución que se obtiene a partir de la de Poisson, en el límite $\mu \gg 1$ ”. Otra es jugar con las distribuciones de números con FDP constante: ¿Qué sucede si los sumamos? ¿Cómo es la FDP de $x_N = \sum_{i=1}^N x_i$, si cada uno de los x_i tiene FDP cte distribuida en (0,1)?



FDP de Gauss

REPASO

Parece haber algún secreto escondido ¿no?



Teorema del Límite Central

Si x_1, x_2, \dots, x_N son variables aleatorias independientes, y cada una de ellas tiene una FDP arbitraria $P_i(x_i)$, con valor medio μ_i y dispersión σ_i^2 entonces

$$x_N = \frac{\sum_{i=1}^N x_i - \sum_{i=1}^N \mu_i}{\sqrt{\sum_{i=1}^N \sigma_i^2}},$$
 se aproxima a una distribución normal para $N \rightarrow \infty$.

$$\lim_{N \rightarrow \infty} P(x_N) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

El caso que mostré es un caso particular de esto, ya que las P_i son siempre la misma PDF, y por lo tanto $\mu_i = \mu = 0.5$ y $\sigma_i = \sigma = 1/\sqrt{12}$.

FDP de Gauss

REPASO

Teorema del Límite Central

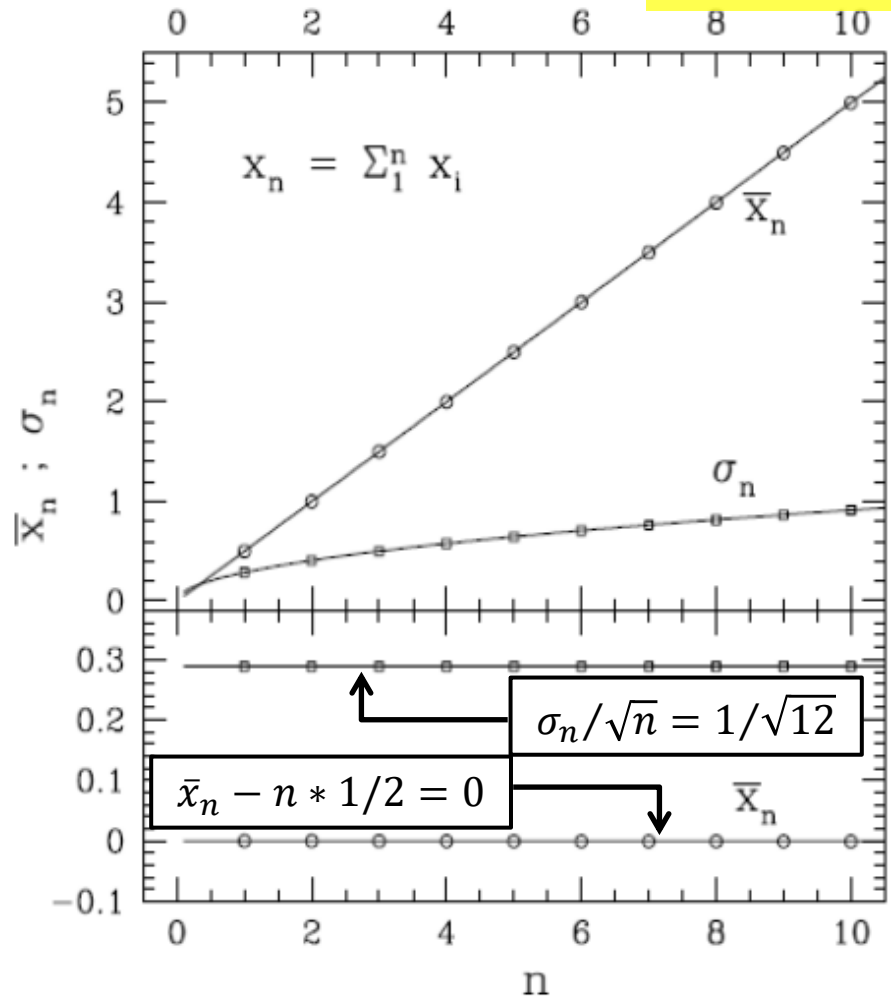
Entonces, para el caso específico de nuestra variable x_N ,

$$\sum_{i=1}^N \mu_i = N * \mu = N * \frac{1}{2}$$

$$\sqrt{\sum_{i=1}^N \sigma^2_i} = \sigma\sqrt{N} = \sqrt{\frac{N}{12}}$$

entonces

$x_N = \frac{\sum_{i=1}^N x_i - N/2}{\sqrt{N/12}}$, se aproxima a una distribución normal para $N \rightarrow \infty$.

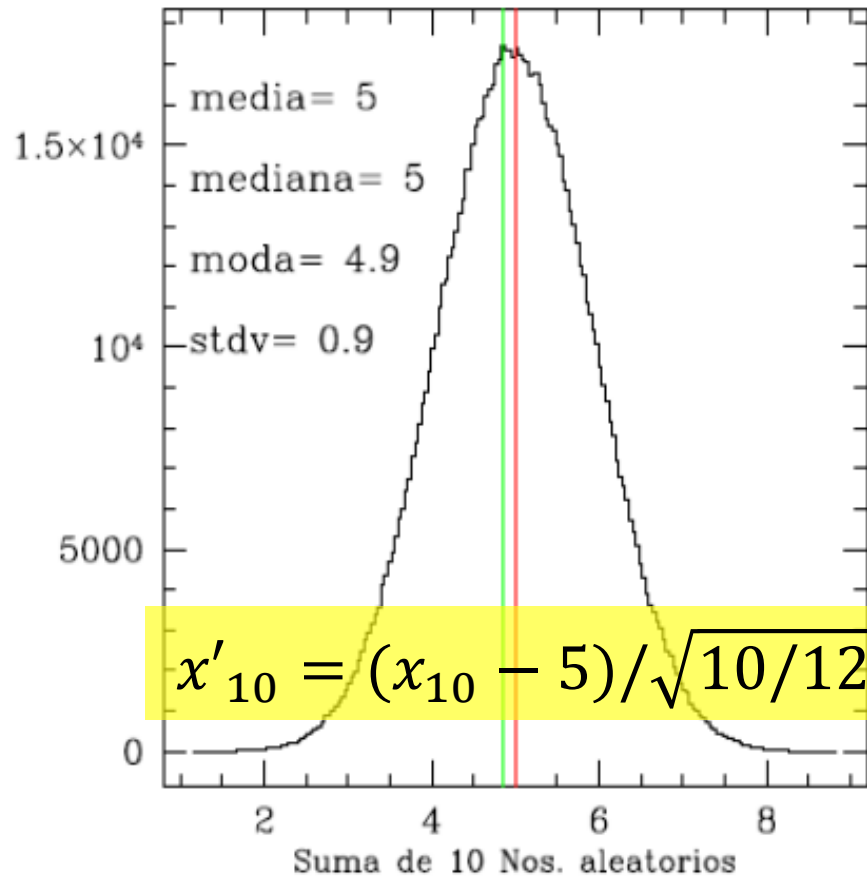


$$\lim_{N \rightarrow \infty} P(x_N) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$$\sum_{i=1}^N x_i \rightarrow \left(\frac{1}{\sqrt{N/12} \sqrt{2\pi}} e^{-\frac{(x-N/2)^2}{2}} \right)$$

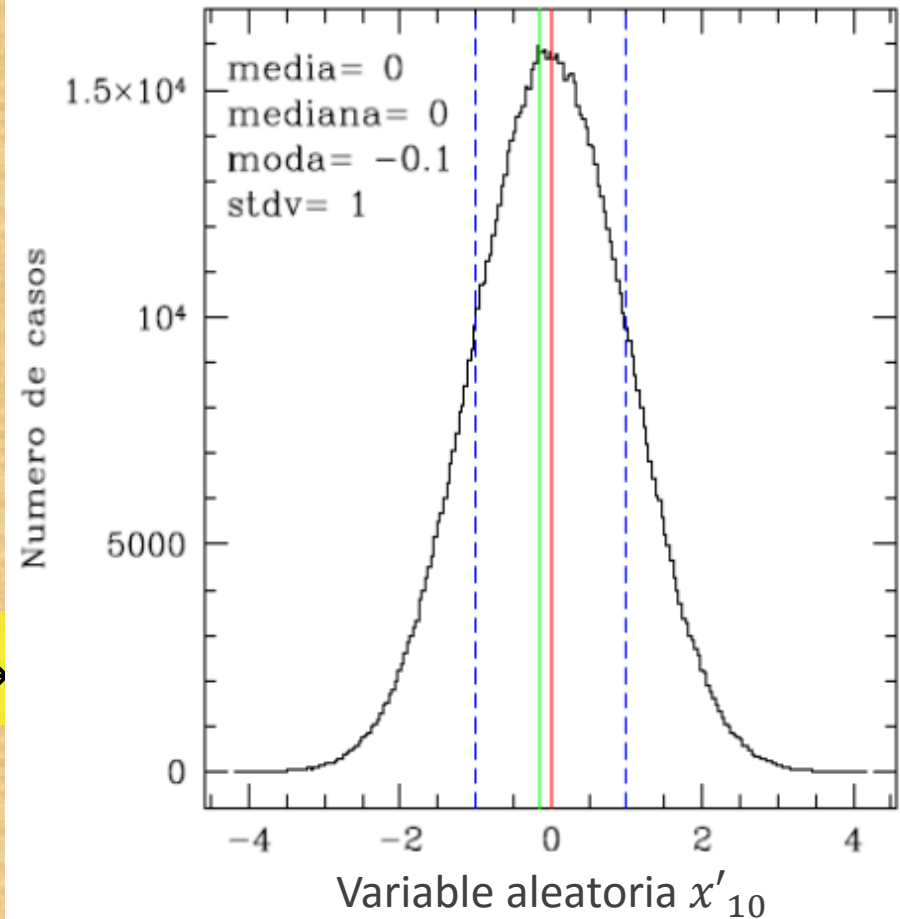
FDP de Gauss

pepo10 ; $N_T = 1000000$; Bin = 0.04



$$x'_{10} = (x_{10} - 5)/\sqrt{10/12} \rightarrow$$

pepo10n ; $N_T = 1000000$; Bin = 0.04

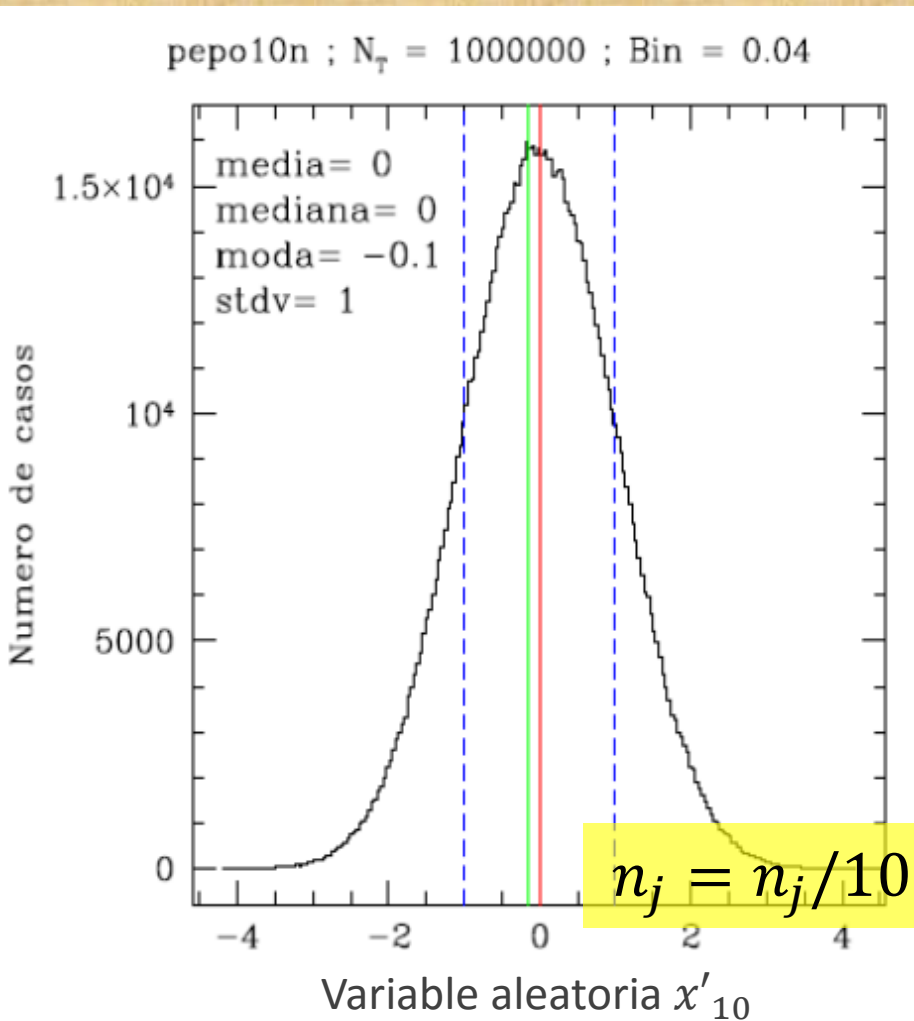


$$P(x_{10}) = N(5, \sqrt{10/12}) = \frac{1}{\sqrt{5\pi/3}} e^{-\frac{(x-5)^2}{2\sqrt{5/6}}}$$

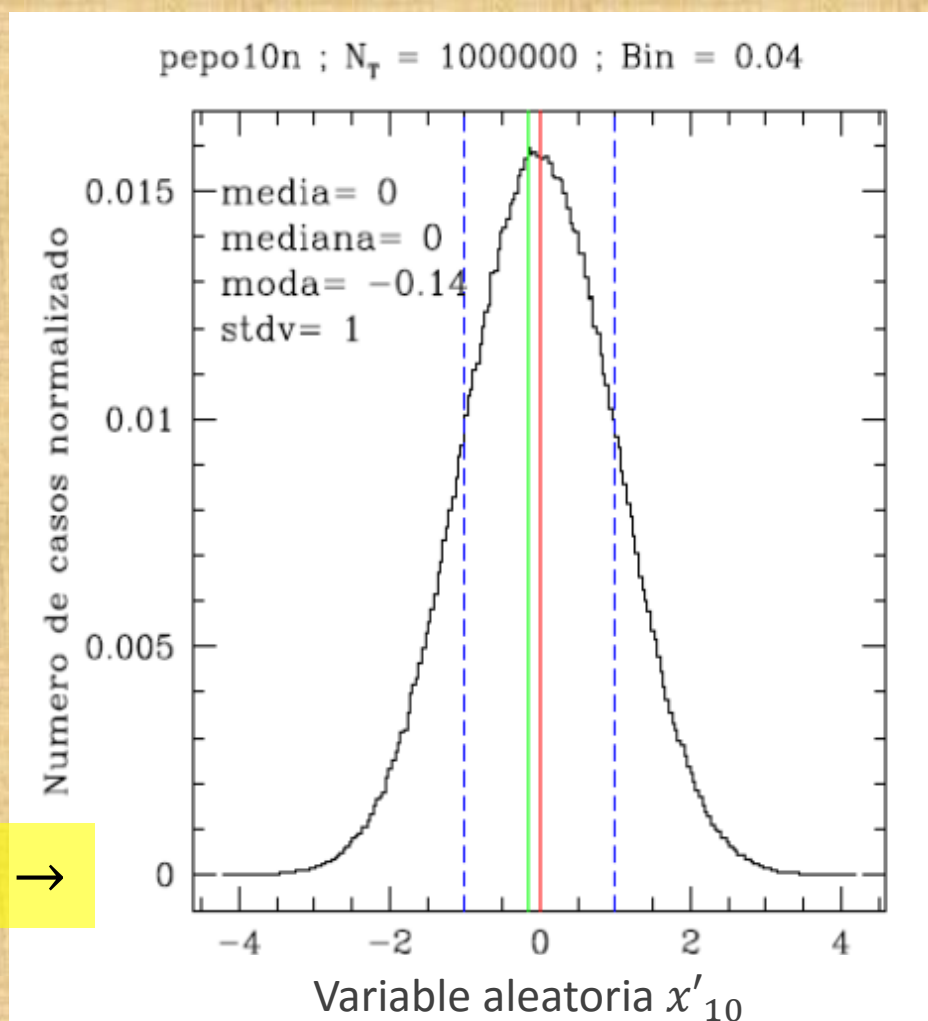
$$P(x'_{10}) = 10^6 N(0,1) = \frac{10^6}{\sqrt{2\pi}} e^{-\frac{x'^2}{2}}$$

FDP de Gauss

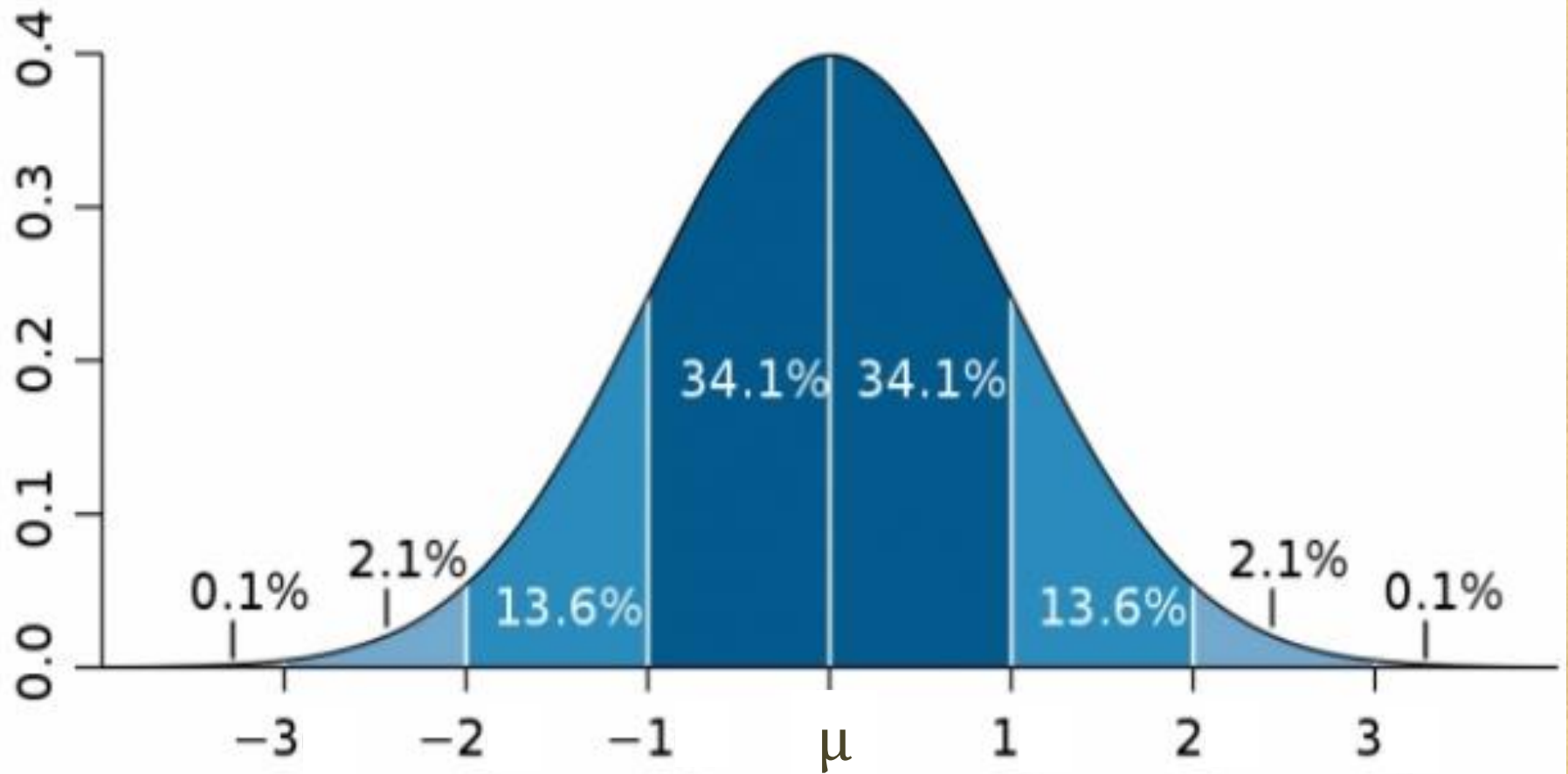
La transformación a un histograma de FDP Gaussiano de $N(0,1)$ todavía no está completa porque el número de casos sigue reflejando la población de 10^6 números. Para corregir esto, dividimos las cuentas del histograma final por 10^6 .



$$n_j = n_j / 10^6 \rightarrow$$

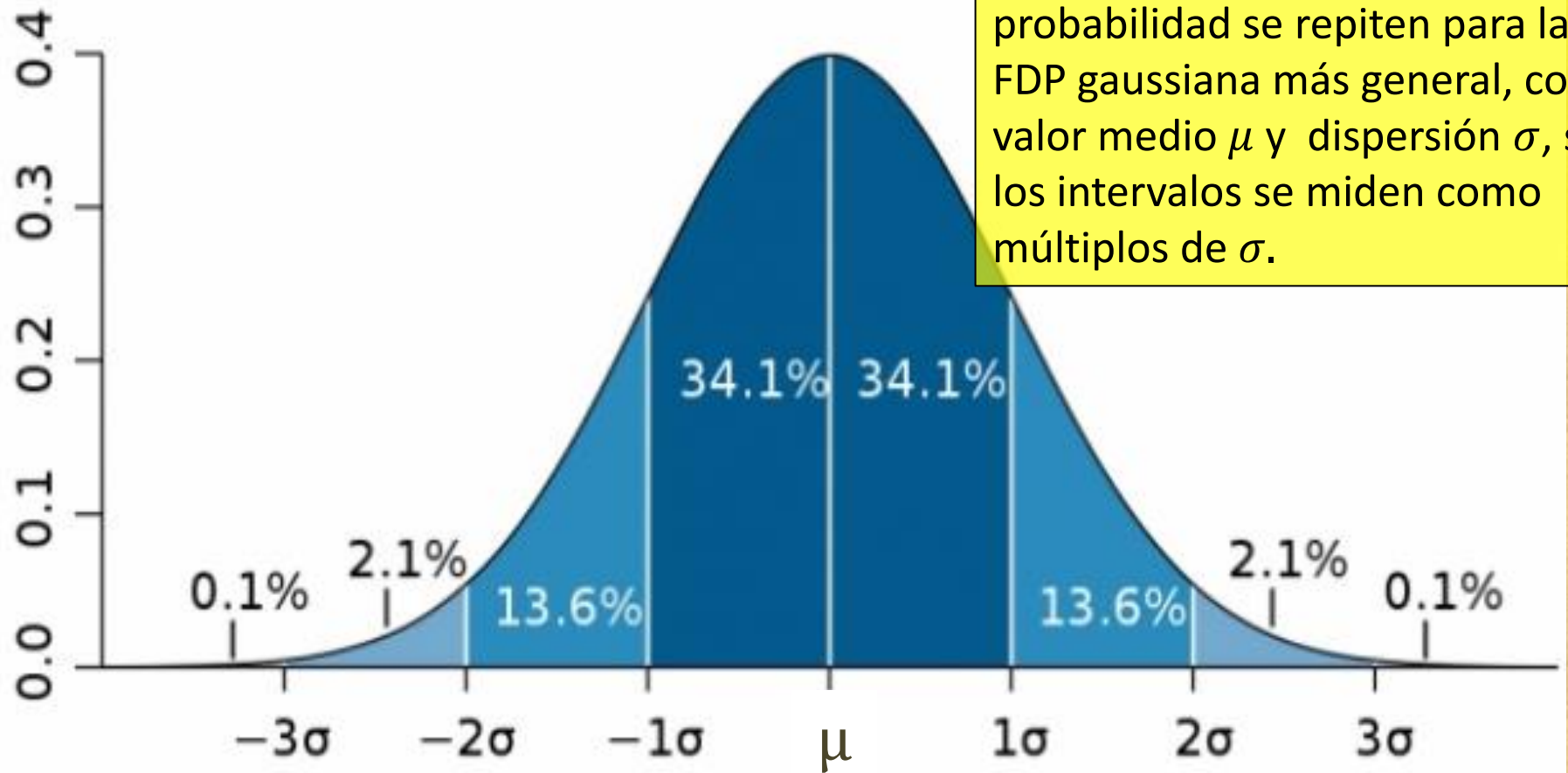


FDP de Gauss



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

FDP de Gauss



Las simetrías e intervalos de probabilidad se repiten para la FDP gaussiana más general, con valor medio μ y dispersión σ , si los intervalos se miden como múltiplos de σ .

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

FDP de Gauss

La forma anterior es la más general de la distribución normal:

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

forma que también está normalizada de forma que su integral en el espacio completo de definición de la probabilidad, $(-\infty, \infty)$ es 1:

$$P_{(-\infty < x < \infty)} = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = 1$$

La FDP de Gauss puede usarse para predecir la probabilidad de que un valor de x esté en un cierto rango de la variable (x_1, x_2) :

$$P_{(x_1 < x < x_2)} = \frac{1}{\sigma\sqrt{2\pi}} \int_{x_1}^{x_2} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \int_{x_1}^{x_2} N(\mu, \sigma) dx$$

FDP de Gauss

La forma usual de llevar a cabo ese cálculo es convirtiendo la $N(\mu, \sigma)$ en $N(0,1)$ con el cambio de variables $x' = (x - \mu)/\sigma$

$$P(x_1 < x < x_2) = \frac{1}{\sqrt{2\pi}} \int_{x'_1}^{x'_2} e^{-\frac{x'^2}{2}} dx' = \text{Erf}(x'_2) - \text{Erf}(x'_1)$$

donde $\text{Erf}(x)$ es

$$\text{Erf}(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{x'^2}{2}} dx'$$

La $\text{erf}(x)$ nos permite calcular un primer “modelo de la realidad” para comparar con nuestras observaciones o con nuestra imaginación.

Si el modelo de la realidad correcto es una distribución $N(\mu, \sigma)$, entonces el valor esperado, $n_{e,j}$, de casos que caerán en el intervalo de ancho Δx centrado en x_j será:

$$n_{e,j} = N_T \{ \text{Erf}(x_j + \Delta x/2) - \text{Erf}(x_j - \Delta x/2) \}$$

donde N_T es el número total de casos (es decir, la suma total del histograma).

FDP de Gauss como modelo de la realidad

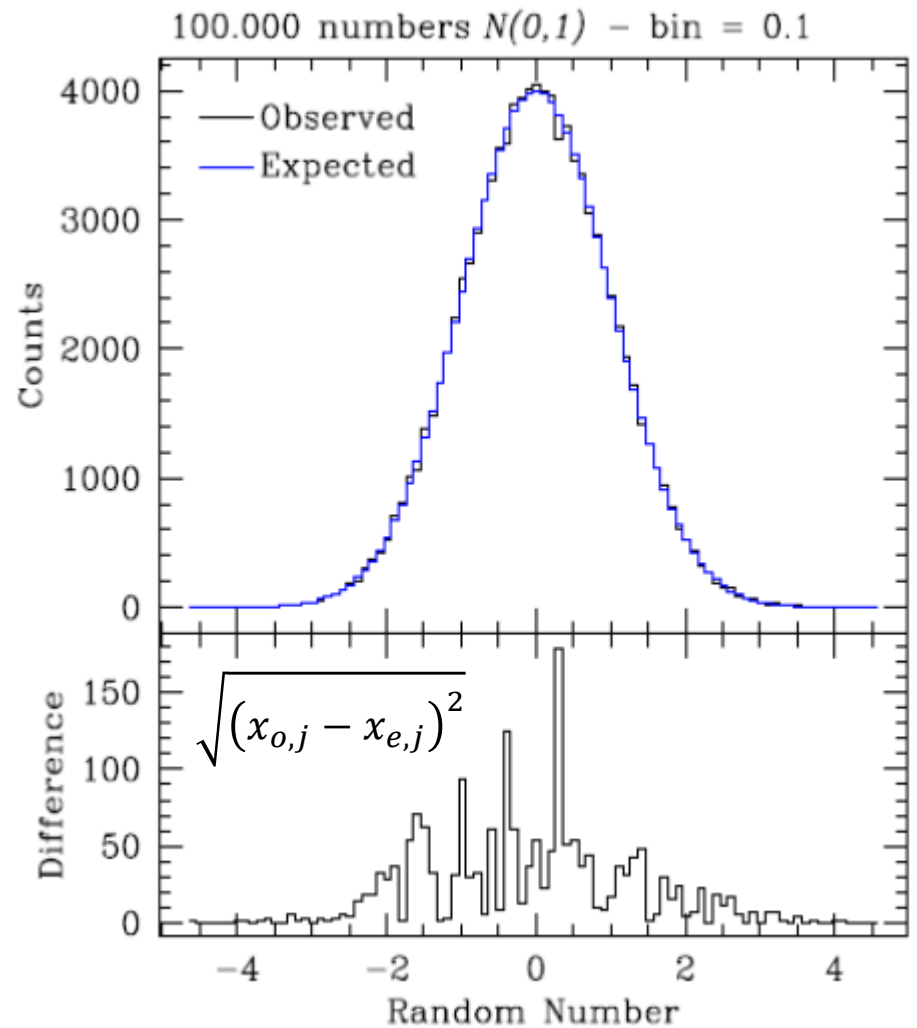
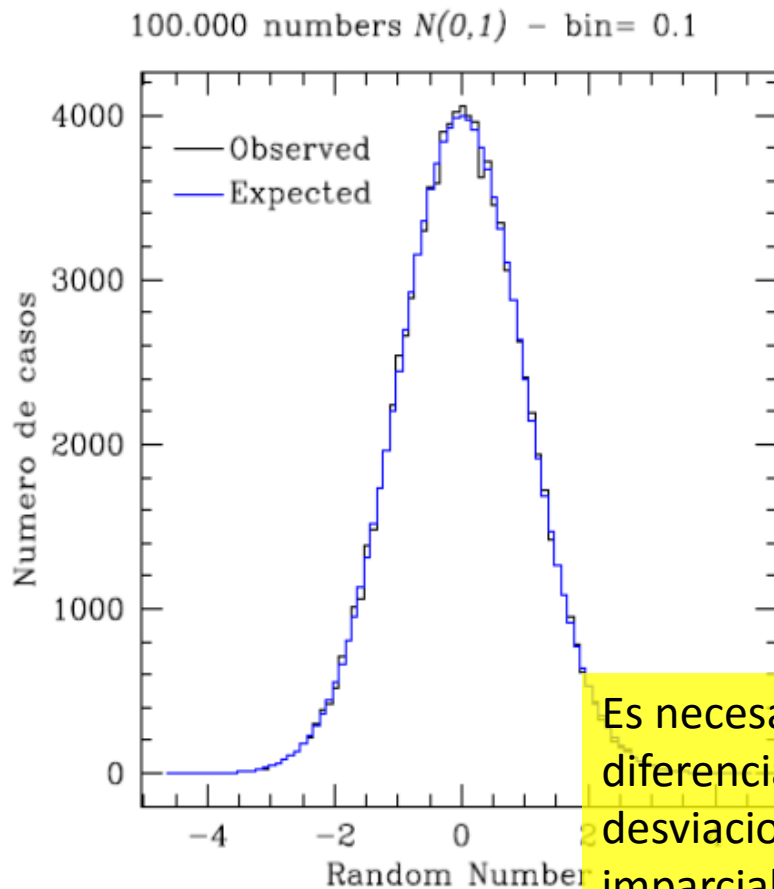
Habiendo hecho esto puedo desarrollar un test cuantitativo para medir cuán diferente es el histograma observado del que predice el modelo. El test que vamos a usar es el llamado χ^2 , que se construye sumando los cuadrados de las diferencias entre las cuentas real, observadas, y las predichas por el modelo, esperadas, bin a bin, normalizadas por las cuentas esperadas. Hay razones fundadas para proponer y usar este estimador, pero por esta clase sólo tomaremos nota de la receta:

$$\chi^2 = \sum_{j=1}^M \frac{(n_{o,j} - n_{e,j})^2}{n_{e,j}}$$

$n_{o,j}$ es el número de casos observados en el bin j , $n_{e,j}$ es el número de casos esperados en el bin j , calculados como se indica en la imagen previa, de acuerdo a la distribución $N(\mu, \sigma)$, donde μ y σ son el valor medio y la dispersión (o desviación estándar), medidas para el histograma que estamos tratando de representar, y la suma se extiende a los M bins que tenga el histograma.

¿Qué es χ^2 ?

$$\chi^2 = \sum_{j=1}^M (n_{o,j} - n_{e,j})^2 ?$$

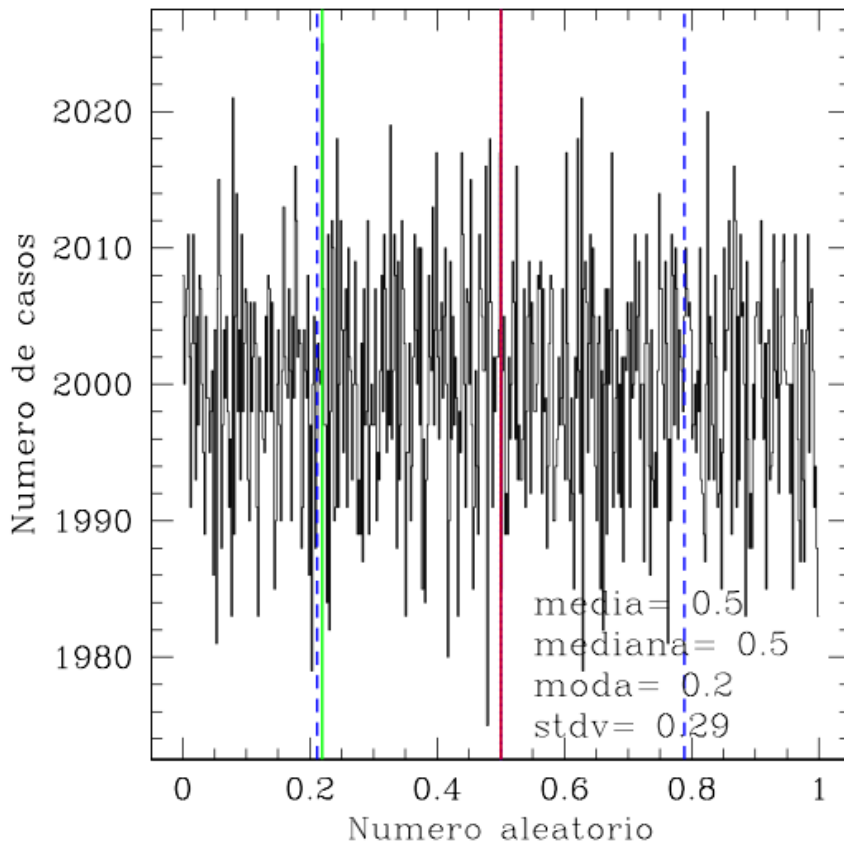


Es necesario normalizar las diferencias para medir las desviaciones de manera imparcial.

$$\chi^2 = \sum_{j=1}^M \left(\frac{n_{o,j} - n_{e,j}}{\sigma_{e,j}} \right)^2$$

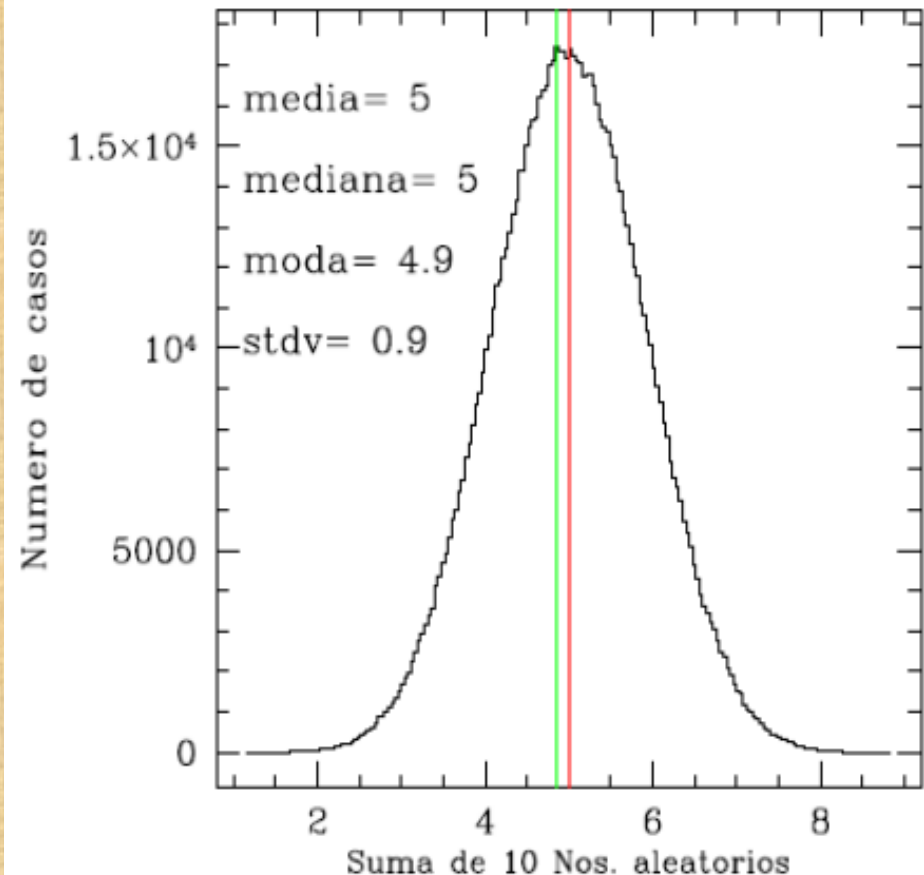
¿Cuál es un buen modelo para σ_e ?

rdn₁e6.dat ; $N_T = 1000000$; Bin = 0.002



El proceso de llenar bins con números sacados al azar una distribución subyacente sigue siempre una FDP de Poisson. μ (y σ) dependerán en general del bin.

pepo10 ; $N_T = 1000000$; Bin = 0.04



El número de cuentas en cada bin sigue una distribución de Poisson, con valor medio μ y $\sigma = \sqrt{\mu}$. En este caso $n_{e,j} = \mu$ (cte.), $\forall j$, pero esto es propiedad de la PDF, no del proceso que llena los bins.

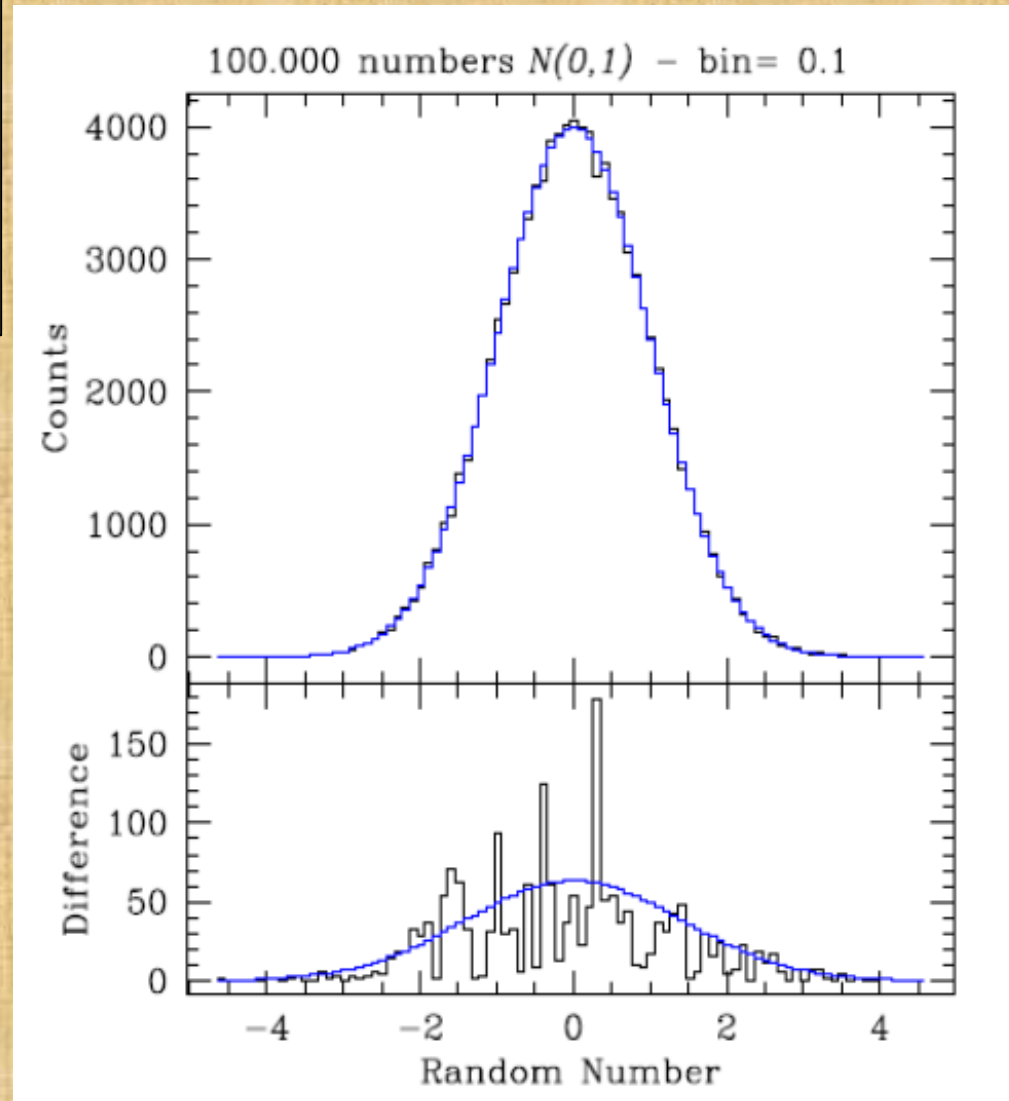
¿Qué es χ^2 ?

En buena medida, esto es lo que ya veíamos en la figura anterior. Las diferencias $|n_o - n_e|$ siguen, aproximadamente la forma $\sqrt{n_e}$. Tenemos entonces:

$$\chi^2 = \sum_{j=1}^M \left(\frac{n_{o,j} - n_{e,j}}{\sigma_{e,j}} \right)^2$$

$$\chi^2 = \sum_{j=1}^M \left(\frac{n_{o,j} - n_{e,j}}{\sqrt{n_{e,j}}} \right)^2$$

$$\chi^2 = \sum_{j=1}^M \frac{(n_{o,j} - n_{e,j})^2}{n_{e,j}}$$

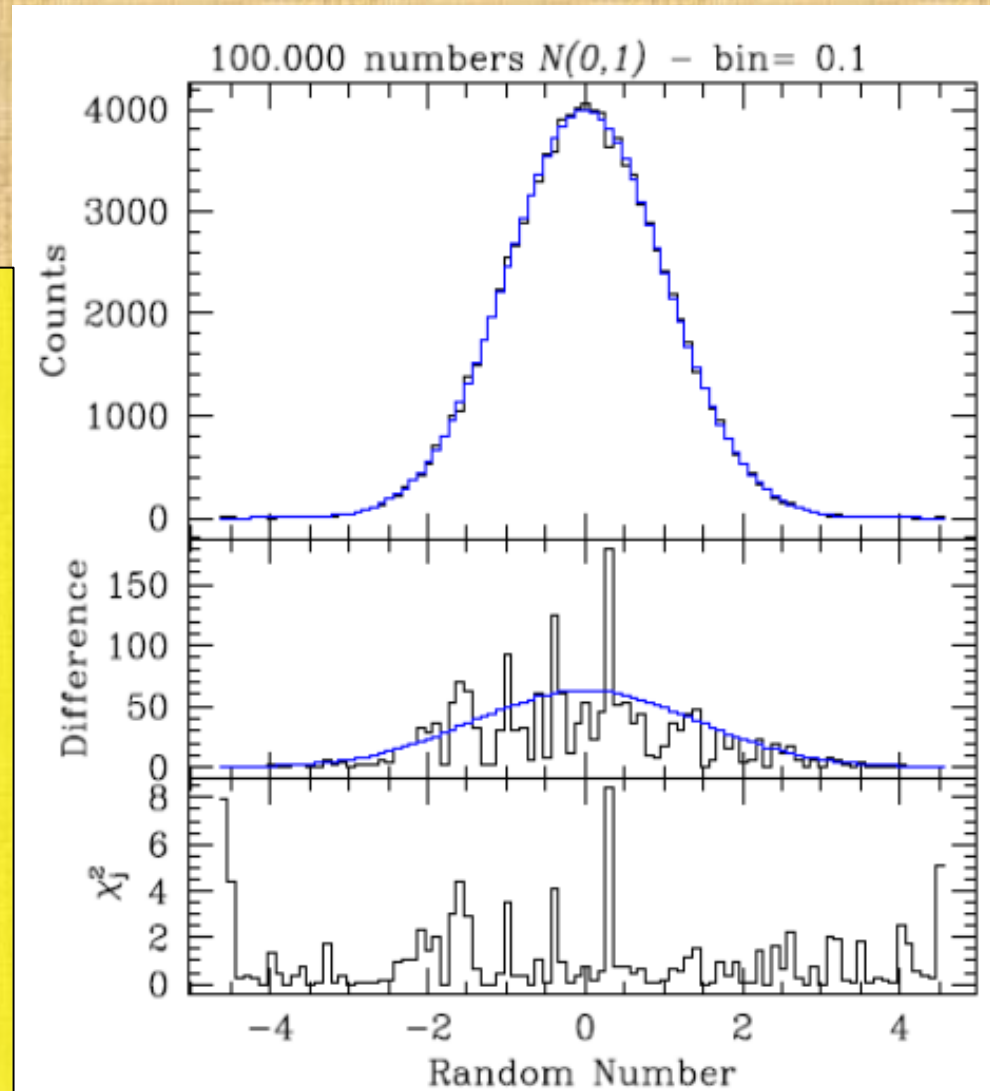


Esta última es la expresión que había puesto como “receta” al final de la clase pasada.

¿Qué podemos esperar de χ^2 ?

$$\chi^2 = \sum_{j=1}^M \left(\frac{n_{o,j} - n_{e,j}}{\sigma_{e,j}} \right)^2 = \sum_{j=1}^M \chi_j^2$$

Esperamos, en principio, que el valor de promedio de $(n_{o,j} - n_{e,j})$ esté bien representado por $\sigma_{e,j}$. Por lo tanto el valor esperado de cada uno de los sumandos del χ^2 es $(n_{o,j} - n_{e,j}) \approx 1$. Si $\chi^2 \gg M$ deberíamos concluir que el modelo de la realidad que codificamos dentro de los $n_{e,j}$ no es una buena representación de los datos. Por otro lado, si tuviéramos $\chi^2 \ll M$ deberíamos concluir que estamos ajustando la realidad por encima de la expectativa estadística (el típico caso de algo “demasiado bueno para ser cierto”).



$$\chi_o^2 = 92.86 ; M = 92$$

Fin de ppt de Clase 5