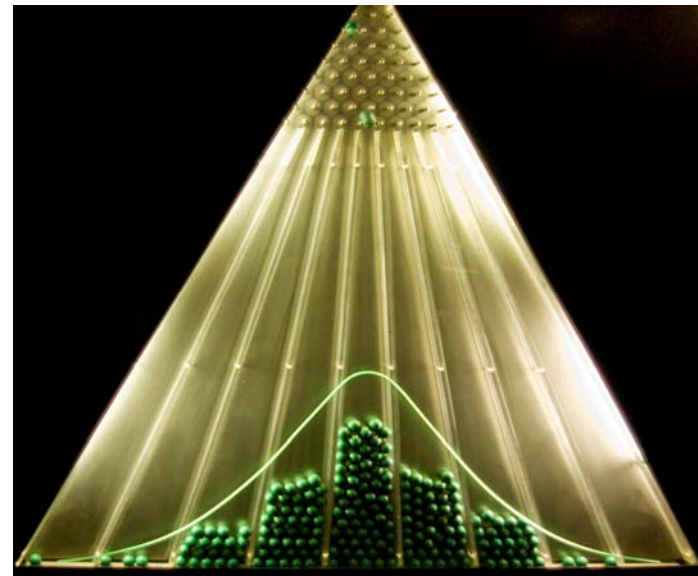


Tema 5 Leyes de los grandes números

1. Muestreo. Inferencia muestral.
2. Ley de los grandes números.
3. Teorema del límite central.
4. Significado del teorema del límite central.
5. Generador de números aleatorios gaussiano.



1. Muestreo. Inferencia muestral

En Física se asocian **funciones densidad de probabilidad** (pdf) a sistemas físicos específicos. La función **pdf** representa el resultado de todas las medidas concebibles sobre dicho sistema si las medidas se pudieran repetir infinitamente bajo las mismas condiciones experimentales.

Una pdf $f(x)$ describe las propiedades de una **población** o **universo**

Población o universo.- Conjunto total de todos los resultados de todas las medidas repetidas infinitamente. Idealización inalcanzable en la práctica.

Espacio muestral.- Conjunto de todos los posibles resultados del experimento. Resultados idénticos se representan por un solo miembro del conjunto.

Muestra (sample).- Conjunto de los resultados de un experimento en particular. Puede haber resultados idénticos.

Cuando realizamos un **experimento**, asociamos un **observable** a una **variable aleatoria** y obtenemos como resultado un conjunto finito de valores o medidas de una cierta cantidad:

$$x_1, x_2, \dots, x_n$$

Se trata de un subconjunto o **muestra** de la **población padre** ("**parent distribution**") o **universo**. Hablamos de una muestra de tamaño n extraída del universo.

1. Muestreo. Inferencia muestral

Supongamos una muestra de tamaño n : x_1, x_2, \dots, x_n . La muestra se puede caracterizar por las siguientes cantidades:

Media muestral

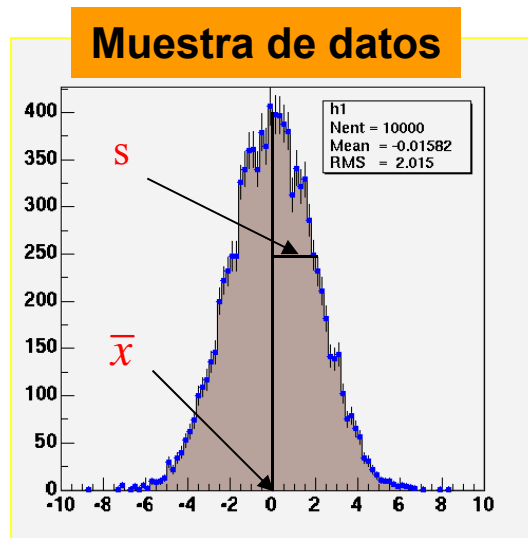
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Varianza muestral

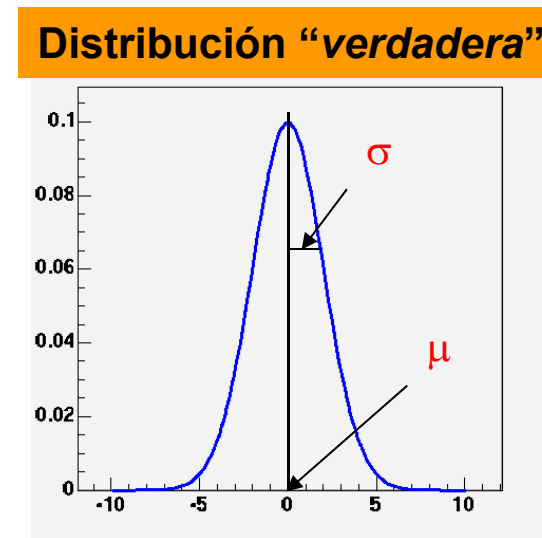
$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Tanto \bar{x}, s^2 son funciones de variables aleatorias (**estadísticos**), por tanto son variables aleatorias también con su propia distribución.

Inferencia muestral Averiguar algo de la “realidad” (distribución **verdadera**) a partir de un número restringido de observaciones. Suponemos que la muestra es representativa de la población.



\bar{x} – media de la **muestra**
 s – Desviación estándar de la **muestra**



μ – media de la dist. **verdadera**
 σ – Desviación estándar de la dist. **verdadera**

1. Muestreo. Inferencia muestral

Al realizar un experimento estamos interesados en la **distribución verdadera** que define nuestro **sistema** o **universo**

Al medir, lo que hacemos es sacar conclusiones ("**inferir**") sobre dicho "**universo**" o "**sistema**" a partir de un número restringido de observaciones

Una medida del valor medio de la población vendrá dada por:

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

\bar{x} es una estimación de la media de la población μ



Una medida de la varianza de la población vendrá dada por:

$$\hat{\sigma}^2 = s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

s^2 es una estimación de la varianza de la población σ^2



El hecho de poner $(n-1)$ en vez de n en el estimador de la varianza se hace para asegurar que se trata de un estimador sin sesgo (Tema 7).

1. Muestreo. Inferencia muestral

Teoremas de Convergencia

Cuando $n \rightarrow \infty$ las propiedades de la muestra se aproximan a las propiedades de la población

$$\hat{\mu} = \bar{x} \rightarrow \mu \quad \hat{\sigma}^2 = s^2 \rightarrow \sigma^2$$

Explicación intuitiva basada en los valores esperados de los estimadores

Valor medio

$$E[\bar{x}] = E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n E[x_i] = \frac{1}{n} \sum_{i=1}^n \mu = \mu$$

$$V[\bar{x}] = V\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n^2} \sum_{i=1}^n V[x_i] = \frac{1}{n^2} \sum_{i=1}^n \sigma^2 = \frac{\sigma^2}{n}$$

La dispersión de $\hat{\mu} = \bar{x}$ en torno a μ disminuye a medida que aumenta n

Varianza

$$E[s^2] = \sigma^2$$

$$V(s^2) = \frac{1}{n}(\mu_4 - \mu_2^2) + \frac{2}{(n-1)n} \mu_2^2$$

La dispersión de $\hat{\sigma}^2 = s^2$ en torno a σ^2 disminuye a medida que aumenta n

2. Ley de los grandes números

Teorema de Chebysev

Sea $h(x)$ una función no negativa de una variable aleatoria x . Siempre se puede encontrar un límite superior a la probabilidad de que $h(x)$ exceda un determinado valor $K > 0$

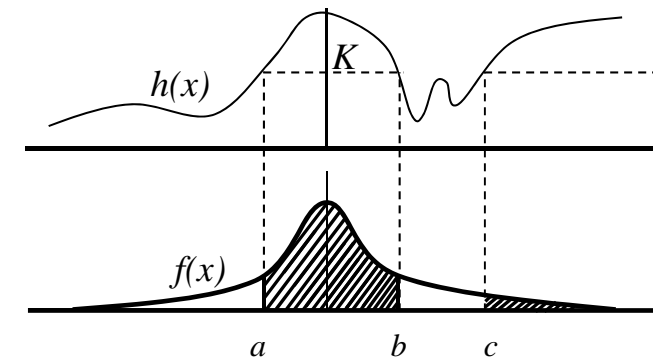
$$\Rightarrow P\{h(x) \geq K\} \leq \frac{E[h(x)]}{K}$$

Demostración

$$E[h(x)] = \int_{\Omega} h(x) f(x) dx \geq 0 \quad \text{Puesto que ambas son funciones no negativas}$$

Definimos R como la región de Ω tal que $h(x) \geq K$

$$\text{Si } x \in R \rightarrow h(x) \geq K \quad \Rightarrow \quad P[h(x) \geq K] = \int_R f(x) dx$$



$$R = [a, b] \cup [c, +\infty)$$

$$E[h(x)] = \int_{\Omega} h(x) f(x) dx \geq \int_R h(x) f(x) dx \geq \int_R K f(x) dx = KP[h(x) \geq K]$$

$R \subset \Omega$

$K \leq h(x) \forall x \in R$

$$P[h(x) \geq K] \leq \frac{E[h(x)]}{K}$$

2. Ley de los grandes números

Desigualdad de Chebysev-Bienayme

Supongamos que x es una variable aleatoria con:

- Valor esperado $E[x]$
- Varianza $V[x]$

$$P\{h(x) \geq K\} \leq \frac{E[h(x)]}{K}$$

$$\begin{array}{ll} h(x) \rightarrow (x - E[x])^2 \geq 0 & \\ K \rightarrow k^2 V(x) & \end{array} \quad \Rightarrow \quad P\left[(x - E[x])^2 \geq k^2 V(x)\right] \leq \frac{E[(x - E[x])^2]}{k^2 V(x)} = \frac{1}{k^2}$$

Tomando raíces cuadradas



$$P\left[|x - E[x]| \geq k V^{1/2}(x)\right] \leq \frac{1}{k^2}$$

Llamando $\mu = E[x]$
 $\sigma^2 = V(x)$



$$P\left[|x - \mu| \geq k\sigma\right] \leq \frac{1}{k^2}$$

k	$1/k^2$	%
1	1	100
2	1/4	25
3	1/9	11
4	1/16	6.3

La probabilidad de que una variable aleatoria x tome un valor fuera del intervalo $[\mu - k\sigma, \mu + k\sigma]$ es menor que $1/k^2$

- Independiente de cual sea la pdf
- Siempre que σ sea conocida.

2. Ley de los grandes números

Ley de los grandes números

Si aplicamos la desigualdad anterior a valores medios

$$\bar{x} = \frac{1}{N} \sum_i x_i$$

$$P[|x - \mu| \geq k\sigma] \leq \frac{1}{k^2}$$

$$P\left[|\bar{x} - \mu| \geq \frac{k\sigma}{\sqrt{N}}\right] \leq \frac{1}{k^2}$$

Mediante el
cambio:

$$\varepsilon = \frac{k\sigma}{\sqrt{N}}$$

$$P[|\bar{x} - \mu| \geq \varepsilon] \leq \frac{\sigma^2}{N\varepsilon^2} = \delta$$

$$x \rightarrow \bar{x}$$

$$E[\bar{x}] = \mu$$

$$V(\bar{x}) = \frac{\sigma^2}{N}$$

Dado un $\varepsilon > 0$ existe un $\delta > 0$ y un N tal que la probabilidad de que \bar{x} difiera de μ en una cantidad mayor que ε es menor que δ para todo $n > N$

Tomando límites

$$\lim_{N \rightarrow \infty} P[|\bar{x} - \mu| \geq \varepsilon] = 0$$

“La media de n variables aleatorias converge a su valor esperado”

2. Ley de los grandes números

Generalización

En el caso más general tendremos medias de funciones de variables aleatorias, por lo que podemos generalizar del siguiente modo:

Supongamos que tenemos n variables aleatorias x_1, x_2, \dots, x_N todas ellas con valor medio μ y varianza σ^2

Supongamos que tenemos una función $g(x)$ de las variables aleatorias de manera que: $\begin{cases} \text{Valor medio} & E[g] \\ \text{Varianza} & V(g) \end{cases}$

Definimos el valor medio G como: $G = \frac{1}{N} \sum_i g(x_i)$

$$E[G] = \frac{1}{N} \sum E[g] = E[g] \quad V(G) = V\left(\frac{1}{N} \sum g(x_i)\right) = \frac{1}{N^2} \sum V(g(x_i)) = \frac{V(g)}{N}$$

$$P\left[|\bar{x} - \mu| \geq \varepsilon\right] \leq \frac{\sigma^2}{N\varepsilon^2} \quad \longrightarrow \quad \left\{ \begin{array}{l} \bar{x} \rightarrow G \\ \mu = E[G] \\ V(G) = \frac{V(g)}{N} \end{array} \right\} \quad \longrightarrow \quad P\left[|G - E[G]| \geq \varepsilon\right] \leq \frac{V(g)}{N\varepsilon^2}$$

Las medias tienden a sus valores esperados

$$G = \frac{1}{N} \sum g(x_i) \rightarrow E[G]$$

2. Ley de los grandes números

Ejemplo. Integración Montecarlo

Queremos calcular una estimación de la integral definida:

$$I = \int_a^b g(x) dx$$

Supongamos que tenemos N números x_i independientes y aleatorios distribuidos uniformemente, en $[a, b]$ según la pdf:

$$f(x) = \frac{1}{(b-a)} \quad a \leq x \leq b$$

El valor esperado de la función $g(x)$ respecto a la pdf $f(x)$ viene dado por:

$$E[g] = \int_a^b g(x) f(x) dx = \frac{1}{(b-a)} \int_a^b g(x) dx = \frac{I}{(b-a)}$$

Si tomamos N números aleatorios x_i y calculamos el valor medio G como:

$$G = \frac{1}{N} \sum g(x_i) \rightarrow E[G] = E[g] = \frac{I}{(b-a)} \quad \longrightarrow \quad \theta = \frac{(b-a)}{N} \sum g(x_i) \rightarrow \int_a^b g(x) dx$$

La ley de los grandes números nos dice que el estimador θ converge a la solución exacta I cuando la muestra aleatoria tiende a infinito

El valor de θ es una estimador de la integral I

La probabilidad de obtener una estimación de la integral que difiera del valor exacto en una cantidad dada, disminuye a medida que aumenta el número de puntos

3. Teorema del límite central

Ley de los grandes números → La media de una variable aleatoria, o de una función de variables aleatorias tiende a su valor esperado cuando $n \rightarrow \infty$

¿Cómo de rápido ?  Teorema del Límite Central

Supongamos que tenemos n variables aleatorias e independientes x_1, x_2, \dots, x_N con: $x_i \rightarrow \mu_i, \sigma_i^2$

Definimos una nueva variable s de la siguiente forma:

$$s = \sum x_i \quad \longrightarrow \quad \mu_s = \sum \mu_i \quad ; \quad \sigma_s^2 = \sum \sigma_i^2$$

¿Cuál es la distribución de la nueva variable s ?

Independientemente de la distribución original de cada x_i y siempre que las variables sean independientes (Véase Tema 3 → Valor esperado y varianza de una combinación lineal de variables aleatorias)

“Teorema del Límite Central (forma condensada)”

“Sean x_1, x_2, \dots, x_N un conjunto de variables aleatorias independientes de manera que cada x_i está distribuida con un valor medio μ_i y una varianza finita σ_i^2 . Entonces la variable

$$z = \frac{\sum_{i=1}^n x_i - \sum_{i=1}^n \mu_i}{\sqrt{\sum_{i=1}^n \sigma_i^2}} \rightarrow N(0,1)$$

tiene como distribución límite ($n \rightarrow \infty$) una distribución normal centrada en 0 y con varianza 1 ”

3. Teorema del límite central

Demostración (I)

Supongamos para simplificar que todas las variables tienen la misma pdf, el mismo valor medio y la misma varianza:

$$\mu_i = \mu ; \quad \sigma_i^2 = \sigma^2 \quad \forall i=1, \dots, n \quad \longrightarrow \quad z = \frac{\sum x_i - \sum \mu_i}{\sqrt{\sum \sigma_i^2}} = \frac{\sum x_i - \sum \mu}{\sqrt{\sum \sigma^2}} = \frac{\sum (x_i - \mu)}{\sigma \sqrt{n}} = \sum \frac{x_i - \mu}{\sigma \sqrt{n}}$$

Cada variable aleatoria x_i tiene su función característica

$$\Phi(t) = E[e^{itx_i}] = 1 + itE[x_i] + \frac{(it)^2}{2!} E[x_i^2] + \dots = 1 + it\mu + \frac{(it)^2}{2!} (\mu^2 + \sigma^2) + \dots$$

Si las variables x_i son independientes también lo son los términos $\frac{x_i - \mu}{\sigma \sqrt{n}}$

La función característica de z será entonces:

$$\begin{aligned} \Phi_z(t) &= E[e^{itz}] = E\left[e^{it \sum \frac{x_i - \mu}{\sigma \sqrt{n}}}\right] = E\left[e^{it \frac{x_1 - \mu}{\sigma \sqrt{n}}} \cdot e^{it \frac{x_2 - \mu}{\sigma \sqrt{n}}} \dots e^{it \frac{x_n - \mu}{\sigma \sqrt{n}}}\right] = \left(E\left[e^{it \frac{x_i - \mu}{\sigma \sqrt{n}}}\right]\right)^n = \\ &= \left(e^{-\frac{it\mu}{\sigma \sqrt{n}}} E\left[e^{i \frac{t}{\sigma \sqrt{n}} x_i}\right]\right)^n = \left(e^{-\frac{it\mu}{\sigma \sqrt{n}}} \Phi\left(\frac{t}{\sigma \sqrt{n}}\right)\right)^n \end{aligned}$$



$$\Phi_z(t) = \left[e^{-\frac{it\mu}{\sigma \sqrt{n}}} \Phi\left(\frac{t}{\sigma \sqrt{n}}\right) \right]^n$$

3. Teorema del límite central

Demostración (II)

Tomando logaritmos:

$$\ln(1+x) = x - \frac{1}{2}x^2 + \dots$$

$$\begin{aligned} \ln \Phi_z(t) &= n \left\{ -\frac{it\mu}{\sigma\sqrt{n}} + \ln \Phi\left(\frac{t}{\sigma\sqrt{n}}\right) \right\} = n \left\{ -\frac{it\mu}{\sigma\sqrt{n}} + \ln \left[1 + \left(\frac{it}{\sigma\sqrt{n}}\right)\mu + \left(\frac{it}{\sigma\sqrt{n}}\right)^2 \frac{(\mu^2 + \sigma^2)}{2!} + \dots \right] \right\} = \\ &= n \left\{ -\cancel{\frac{it\mu}{\sigma\sqrt{n}}} + \cancel{\frac{it\mu}{\sigma\sqrt{n}}} + \frac{1}{2!} \left(\frac{it}{\sigma\sqrt{n}}\right)^2 (\cancel{\mu^2} + \sigma^2) + \dots - \frac{1}{2} \left(\frac{it}{\sigma\sqrt{n}}\right)^2 \mu^2 + \dots \right\} = -\frac{t^2}{2} + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) \end{aligned}$$

$$\ln \Phi_z(t) = -\frac{t^2}{2} + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) \xrightarrow[n \rightarrow \infty]{} \Phi_z(t) = e^{-\frac{t^2}{2}}$$

Función característica de una distribución normal con media 0 y varianza unidad

En el límite de grandes números z se distribuye según una gaussiana $N(0,1)$

4. Significado teorema del límite central

¿Por qué la mayoría de las medidas que hacemos responden a distribuciones gaussianas?

La mayoría de las cosas que medimos dependen de multitud de efectos y factores perturbadores e incontrolables. Lo que medimos es lo que queremos medir ("**efecto verdadero**") más un gran número de pequeñas contribuciones ("**error total**"):

Medida = "Efecto verdadero" + "Error total"



$$\text{"Error Total"} = \sum_i (\text{factores})_i$$



Cada factor responde a una distribución o pdf desconocida.

Combinación de un gran número de errores (variables aleatorias) independientes

*Cuando $n \rightarrow \infty$ el error total se distribuye **normalmente***

¿Cómo se distribuyen los valores medios?

Según el TLC, los valores medios se distribuyen de acuerdo a una distribución normal con valor medio, el valor medio original μ y desviación estándar:

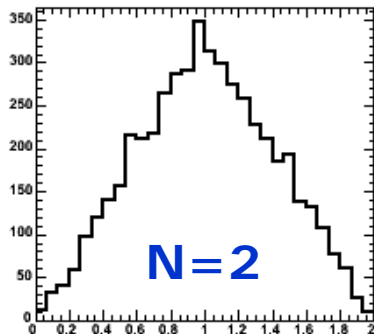


$$\frac{\sigma}{\sqrt{n}}$$

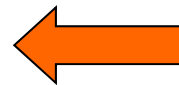
¡ El error de la media se reduce a medida que realizamos más medidas en la media !

5. Generador de números aleatorio gaussiano

5000 números aleatorios distribuidos uniformemente entre [0,1]
Media $\frac{1}{2}$, Varianza $\frac{1}{12}$



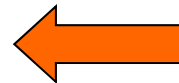
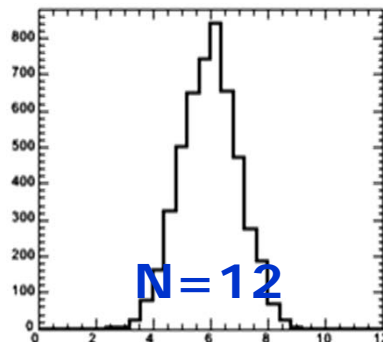
5000 números aleatorios cada uno de los cuales es la suma de dos números aleatorios uniformes distribuidos entre [0,1]



$$x = x_1 + x_2$$

5000 números aleatorios, suma de tres números aleatorios uniformes distribuidos entre [0,1]

$$x = x_1 + x_2 + x_3$$



5000 números aleatorios, suma de 12 números aleatorios uniformes distribuidos entre [0,1]

