

ANÁLISIS ESTADÍSTICO

MUESTREO DE POBLACIONES FINITAS

Jorge Fallas
jfallas56@gmail.com

2010

Temario

- **¿Porqué muestrear?**
- **¿Para qué muestrear?**
 - Estimar parámetros de población
 - Prueba de hipótesis
 - Exploratorio
- **Poblaciones, muestras y estimadores**
 - Infinitas
 - Finitas
 - Propiedades de estimadores
- **Esquemas de muestreo**
 - Probabilístico
 - Al azar
 - Simple al azar
 - Diseños más complejos
 - No Probabilístico
 - Sistemático
 - Experto
 - Conveniencia
- **Herramientas**
 - XLStats
 - Instat

¿Porqué muestrear?

- No es posible medir la totalidad de los elementos
 - Dinero, tiempo, herramientas
 - No es necesario
- Existen métodos estadísticos que nos permiten realizar inferencias a partir de muestras
- Deducción: Investigador(a) no posee datos sobre la totalidad de la población y lo que hace es generalizar la conclusión obtenida a partir de una muestra a aquellos individuos no medidos o estudiados

Población y muestra

- **Población: DE elementos, De observaciones**
- **Infinita** (Elementos no se pueden contar)
 - Árboles de Mundo
 - Estrellas del Universo
 - Tiempo de respuesta a un estímulo
- **Finita** (Elementos se pueden contar)
 - Árboles en campus UNED
 - No. estudiantes curso estadística
 - Venados cola blanca en San Lucas
 - Pobladores de San José de la Montaña

PARÁMETROS Y ESTIMADORES

- **Parámetro** es una propiedad de la población de observaciones (Ej. media, desviación estándar). Es una constante desconocida. Representado con letras griegas: α , β , τ y σ
- **Estimador** es una función que permite estimar el valor del parámetro a partir de los datos de una muestra. Representado con letras latinas: a , b , r y s .
- El estimador es una variable aleatoria- Su valor cambia de muestra en muestra y por tanto posee una distribución de probabilidades conocida como Distribución muestral del estimador
 - Insesgado
 - Variancia mínima
 - Consistente
- **Propiedades**
 - Insesgado
 - Eficiente
 - Suficiente
 - Consistente



Estimadores: Propiedades

- **Insesgado**: valor esperado = parámetro
 - Sesgo = $E(\hat{\theta}) - \theta$ = Diferencia entre el valor esperado del estimador y el valor del parámetro

- **Eficiencia** – Mínimo Cuadrado medio del error

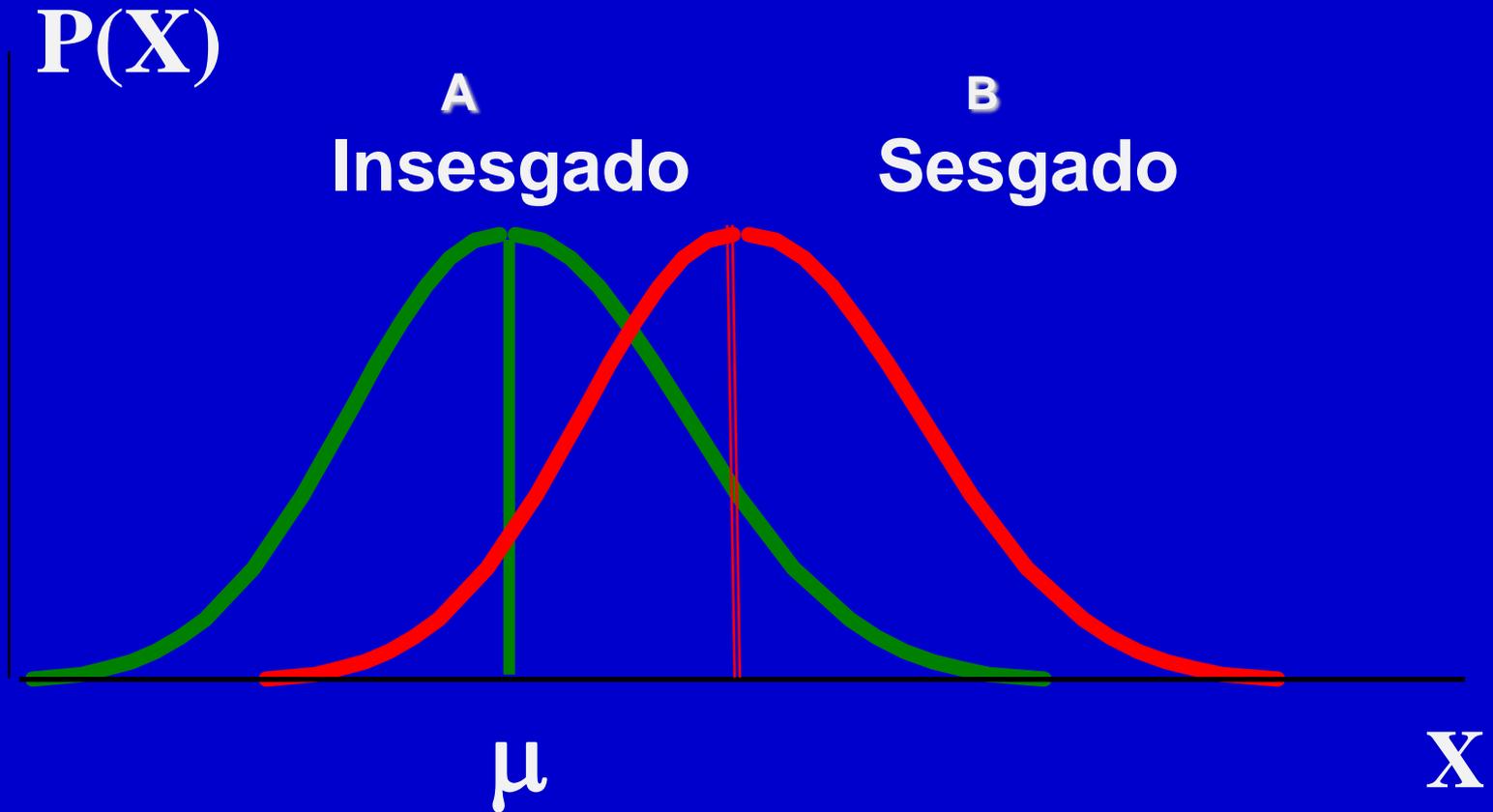
$$\text{M.S.E.} = E(\hat{\theta} - \theta)^2$$

$$= E[\hat{\theta} - E(\hat{\theta})]^2 + [E(\hat{\theta}) - \theta]^2$$

(varianza + (sesgo)²)

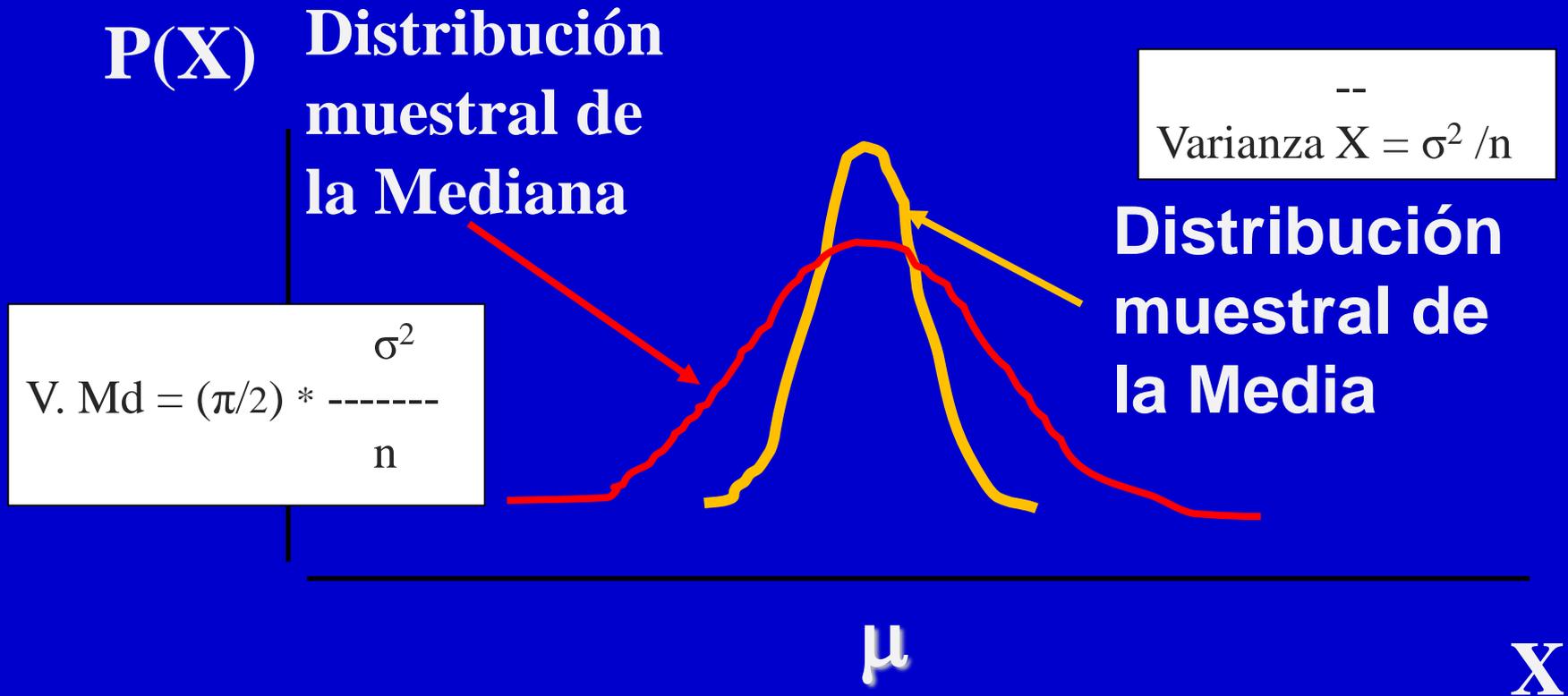
- Que tan bueno es el estimador cuando hacemos predicciones.
- Se prefiere aquel estimador que tenga el error medio cuadrático más pequeño alrededor del parámetro

Insesgado



¿Porqué el estimador en A es insesgado y en B es sesgado?

Eficiencia de mediana y media



Eficiencia – Mínimo Cuadrado medio del error

Propiedades estimadores

- **Consistencia**

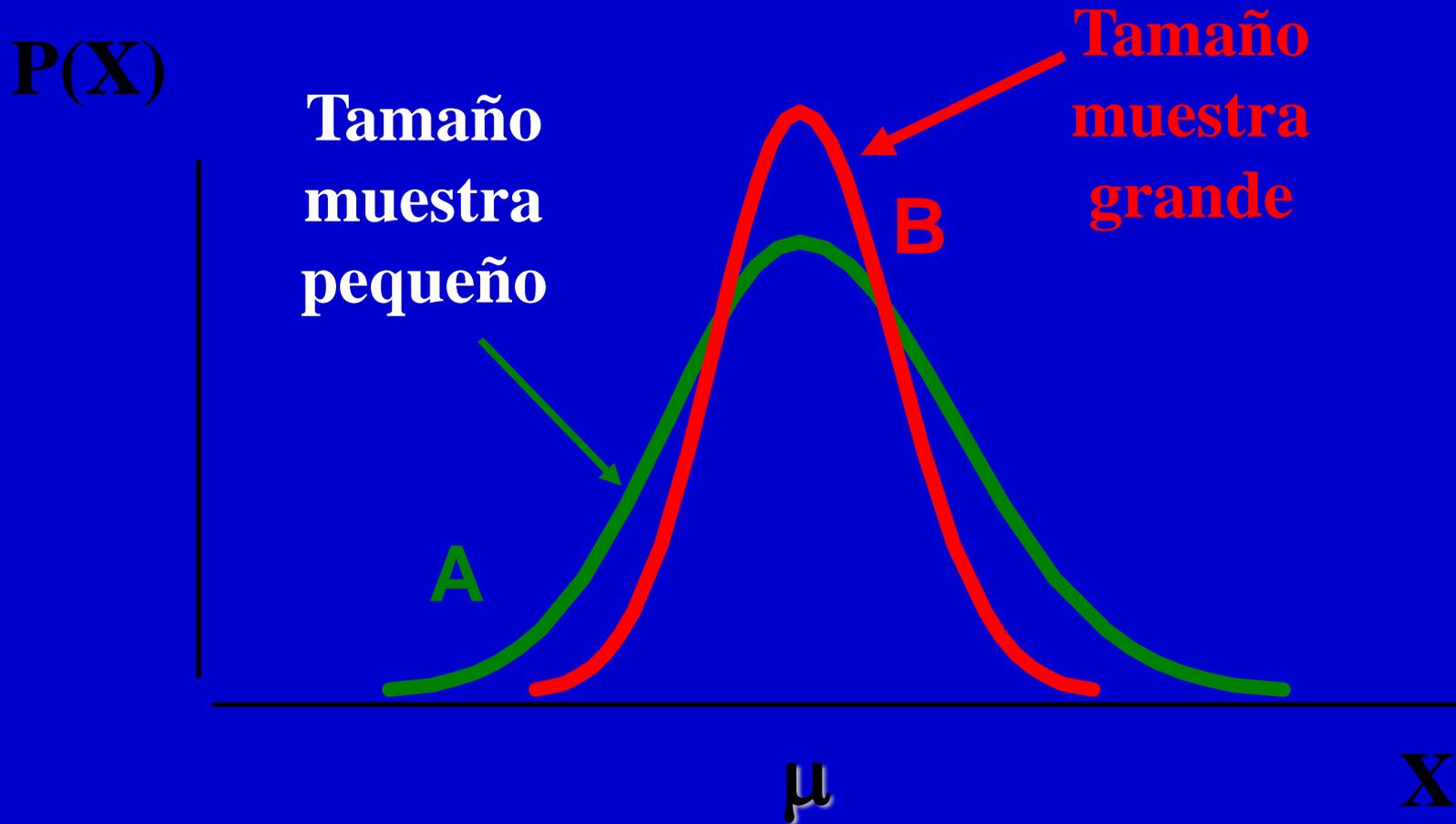
$$\lim_{n \rightarrow \infty} \hat{\theta} = \theta$$

- Al aumentar el tamaño de la muestra, la variación del estimador con respecto al valor del parámetro se reduce

- **Suficiente**

- Todos los datos son utilizados en la estimación

Consistencia



Estimaciones de μ
Media, Mediana, Rango SIC

Insesgado
Media
Mediana

Sesgado
Rango SIC

Inconsistente
Rango SIC

Consistente
Media
Mediana

Varianza mínima
Media

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

Error estándar:
al aumentar n
decrece el error

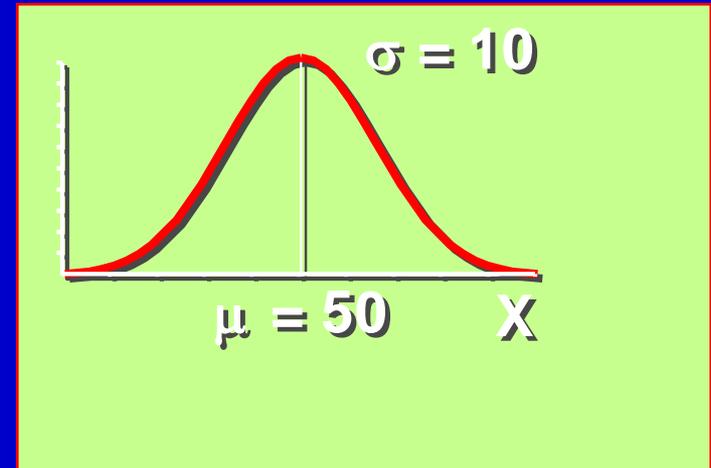
Media aritmética: estimador insesgado,
de varianza mínima y consistente

Población normal → muestra normal

Distribución Población

Tendencia Central

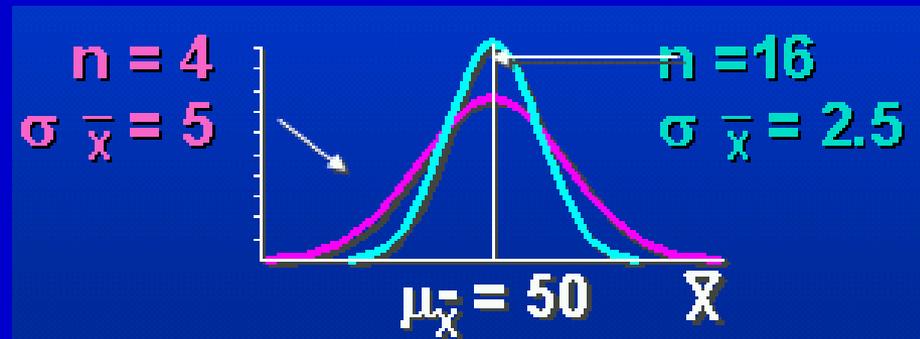
$$\mu_{\bar{X}} = \mu$$



Variación

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

Distribuciones muestrales



Teorema del límite central

- Conforme el tamaño de muestra incrementa la distribución muestral de los promedios se aproxima una distribución normal con media μ y varianza σ^2/n

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

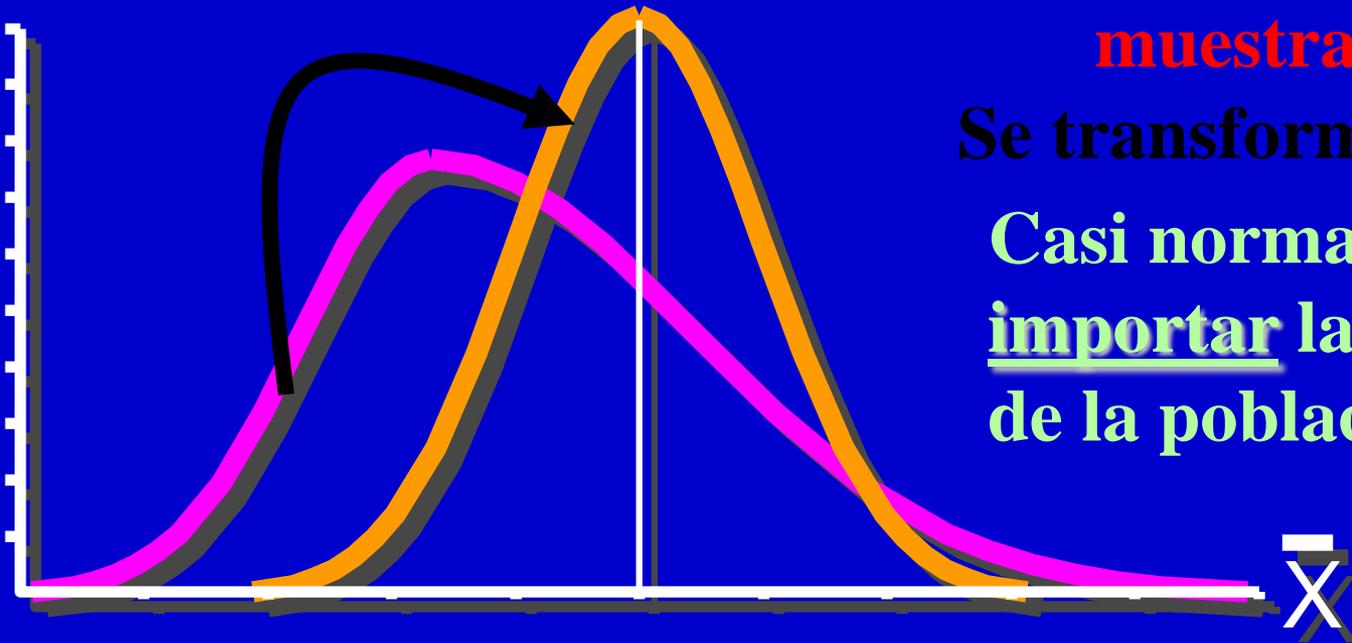
Teorema del límite central

Dado un tamaño de muestra
suficientemente grande

**La distribución
muestral**

Se transforma en

**Casi normal sin
importar la forma
de la población**



Pruebas de normalidad

- **W de Shapiro-Wilk**
- Considerada una de las pruebas mas fidedignas o imparciales.
- **D de Kolmogorov-Smirnov**
- Prueba de bondad de ajuste de la distribución normal
- **Lillifors**
- Versión modificada de la D de Kolmogorov-Smirnov

- **H₀**: los datos provienen de una distribución normal
- **H_a**: los datos no provienen de una distribución normal
- Ver Gráfico de probabilidad normal

XLStat

Options for XLStatistics [X]

Entering Data

If you have a range of cells highlighted in a workbook when an XLStatistics analysis workbook is selected from the XLStatistics menu, you can choose to automatically enter it into the workbook's Data area.

Don't automatically enter data

Automatically enter data

Protection of workbooks and worksheets

The XLStatistics workbooks are templates that can be corrupted by certain operations. Protecting the workbooks stops accidental changes to the structure of the workbook (like renaming sheets); protecting worksheets allows only safe changes to cells (those with blue text). Unfortunately protecting worksheets also stops you from formatting - if necessary, you can unprotect individual sheets by using Tools | Protection | Unprotect sheet - no password is set.

Open XLStatistics workbooks with Protection of workbooks turned on

Open XLStatistics workbooks with Protection of worksheets turned on

Recording

You can use Record on the XLStatistics menu to save results in workbooks with linked data.

Use Ctrl-r as shortcut key for recording

Show Record sheet after recording

Hide Launchpad

Open XLStatistics workbooks in Read-Only Mode

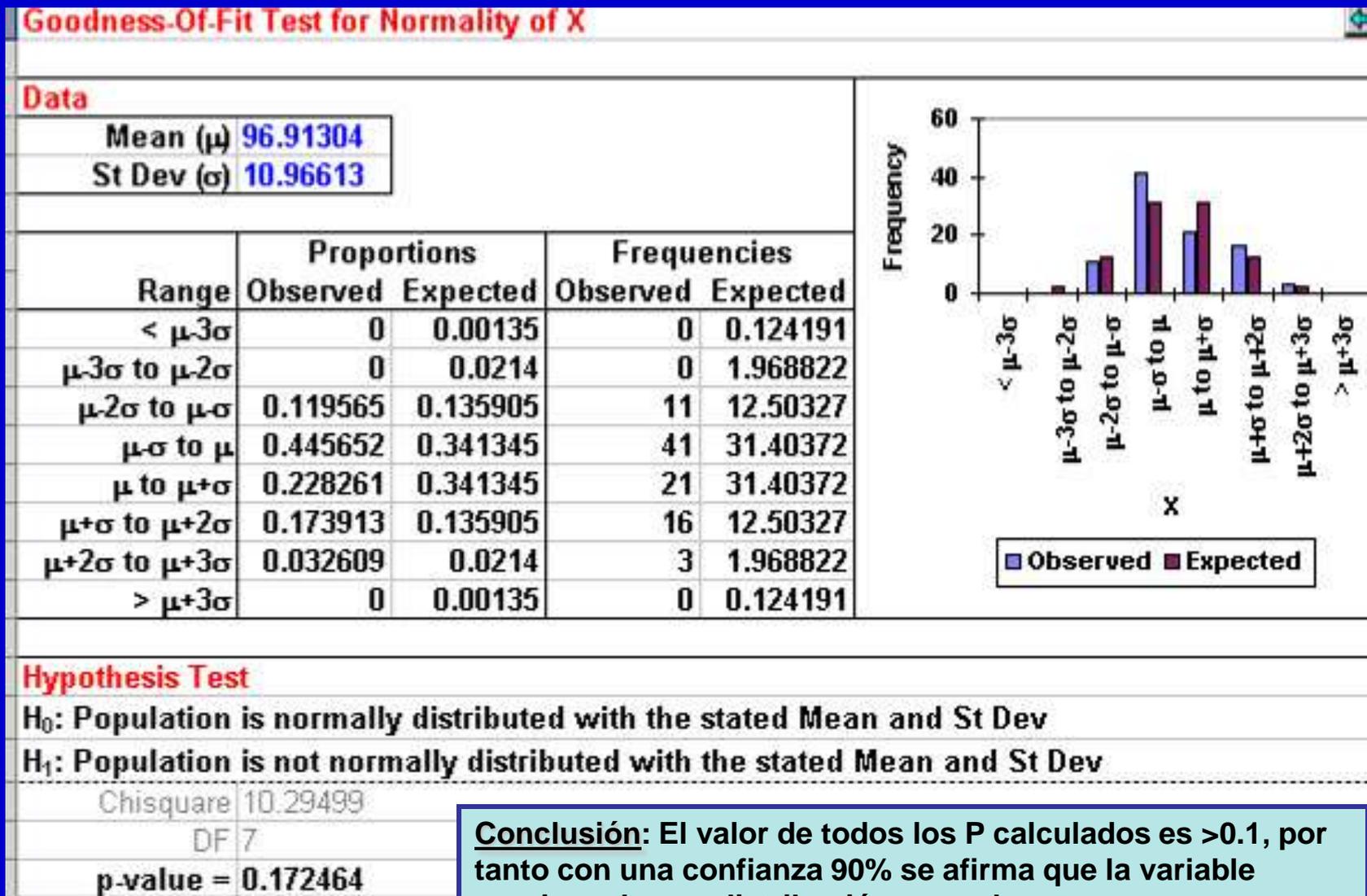
Save current options

Zoom sheets to fit current windows on opening

Display 'Paste Special' message when copying charts

Para carga el macro haga un doble clic sobre el archivo XLSTAT.xls. XLSTAT es un conjunto de funciones estadísticas y matemáticas que le permiten visualizar y analizar sus datos. Utiliza como interface de entrada y salida las hojas de cálculo Excel.

XLStats: prueba normalidad



Conclusión: El valor de todos los P calculados es >0.1 , por tanto con una confianza 90% se afirma que la variable proviene de una distribución normal

Muestreo probabilístico



- Muestreo probabilístico parte del supuesto de que se conoce cada una de las unidades de muestreo que componen la población y que por lo tanto se pueden listar **todas las posibles muestras**, las cuales se conocen como **espacio o marco muestral**
- Permite evaluar el sesgo y el error estándar de los estimadores así como la eficiencia de diferentes diseños de muestreo.
 - Permite estudiar ventajas, desventajas, limitaciones y alcances de dos o más propuestas de muestreo y seleccionar aquella que mejor se ajuste a las características de la población en estudio

Muestreo al azar y simple al azar

- **Muestra al azar**
 - Muestreo con reemplazo
 - Muestras se obtienen en forma **independiente** de una **población** de unidades **infinitas**. En este esquema de muestreo cada elemento tiene la misma probabilidad de ser seleccionada ($P(A) = n/N$).
- **Muestra simple al azar**
 - Muestreo sin reemplazo
 - Muestras se obtienen en forma **no independiente** de una **población** de unidades **finitas**. En este esquema de muestreo cada elemento tiene probabilidad $P(A/B) = (n-1)/(N-1)$ de ser seleccionada.
- Muestreo con reemplazo es *ineficiente*; ya que el seleccionar un individuo más de una vez no aporta información adicional sobre la población.

Intensidad muestreo

- Dado por n/N en donde n es el número de muestras seleccionadas y N el total de muestras de la población.
- Cuántas muestras debo seleccionar ? ó ¿Qué tan grande debe ser la muestra?

Estimación tamaño muestra

$$n = \frac{t^2 * S^2 * N}{N E^2 + t^2 * S^2}$$

S^2 se puede cambiar por CV%

E: error permisible

S: desviación estándar

t: valor de t de Estudiante para $(1 - \alpha/2)$ con $n-1$ grados de libertad; $n > 30$ utilizar Z ó 2 como aproximación.

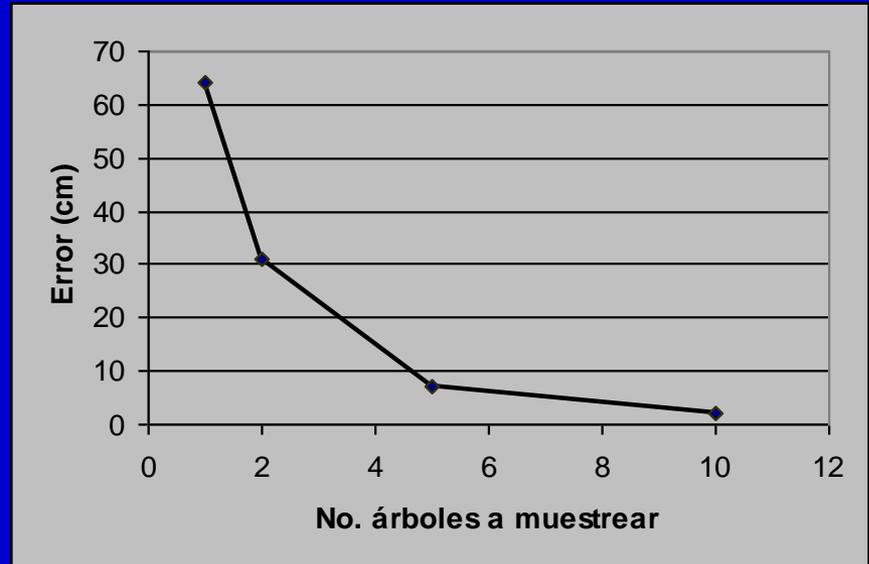
N: Número de muestras o elementos en la población.

Ejemplo

Datos:

$E = 2\text{cm}$ $N = 100$ $t = 2$

Error (cm)	Tamaño de Muestra
± 1	64
± 2	31
± 5	7
± 10	2



1. Al menor error mayor tamaño de muestra
2. A menor variabilidad menor tamaño de muestra para un error dado

CONCLUSION: TAMAÑO DE MUESTRA DEPENDE DE Variabilidad población, Error deseado y nivel de confianza

XLStatistics

XLStatistics - Excel Workbooks for Statistical Data Analysis

© Rodney Carr 1997-2003

Data Analysis Workbooks

Number of Variables	1	1 Numerical 1Num	1 Categorical 1Cat	
	2	1 Numerical 1 Categorical 1Num1Cat	2 Categorical 2Cat	2 Numerical 2Num
	3	1 Numerical 2 Categorical 1Num2Cat	2 Numerical 1 Categorical 2Num1Cat	
	n	1 Numerical n Categorical 1NumnCat	n Categorical nCat	n Numerical nNum
		n Numerical 1 Categorical nNum1Cat		

Other Workbooks

Probability Functions PDF	Transform Transform	Options	Help
Sample Selection SampSel	Populate Populate	<input type="checkbox"/> Hide Launchpad	<input type="checkbox"/> Zoom sheets to fit window
Quality Control Control			

Menú general

Muestras al azar: XLStats

Sample Selection

SampSel

EZRStats ? Ad

Analysis ▶



Univariate Statistics

	A	B	C	D	E	F
1	Taking Random Samples from a Population					
2						
3	Population		Sample Size	10		Sample
4		16.7				10.5
5		28.4				27.4
6		19.5				9.5
7		19.9				15.5
8		17.6				9.9
9		16				15.1
10		18.2				12.5
11		15.6				16
12		14.8				20.5
13		16.1				17.7
14		31.7				
15		10.9				
16		14.4				
17		15.7				
18		24.7				
19		21.9				
20		6.9				
21		24.2				
22		37.9				
23		22.1				
24		12				

Sample Size 10

Sampling method
 With replacement
 Without replacement

New Sample

Univariate Statistics

Select range and then click Process a range with numbers in order to obtain Univariate Statistics

Average: 15.4600
Std Dev: 0.2338
Var: 30.3560
Skewness: 0.0016
N: 10
min: 9.5
max: 27.4
range: 17.9
missing: 0
Negative: 0
Positive: 10
Zero: 0
10th percentile is 9.5
25th percentile is 9.9
50th percentile is 15.1

Process

Finished / Cancel

Muestreo Estratificado

- Población a muestrear es muy heterogénea pero se puede dividir en estratos homogéneos a su interior pero heterogéneos entre sí.
- Se obtiene una nuestra simple al azar de cada estrato. La muestra agregada se denomina una *muestra aleatoria estratificada.*
- Intensidad de muestreo puede ser diferente para cada estrato y está dado por n_h / N_h ; en donde n_h es el número de elementos seleccionados en el estrato h y N_h es el total de elementos del estrato h .

Ventajas

- Permite utilizar información sobre una o más características de la población para delimitar estratos homogéneos a su interior.
- El muestreo al azar estratificado es más preciso que el muestreo simple al azar siempre y cuando la variabilidad al interior del estrato sea menor comparada con la variabilidad entre estratos.
- El **error de muestreo** está dado por la variabilidad entre unidades de muestreo de cada estrato y no por la variabilidad entre estratos. La ganancia en precisión se manifiesta en una reducción del error de muestreo

Estratos: coberturas

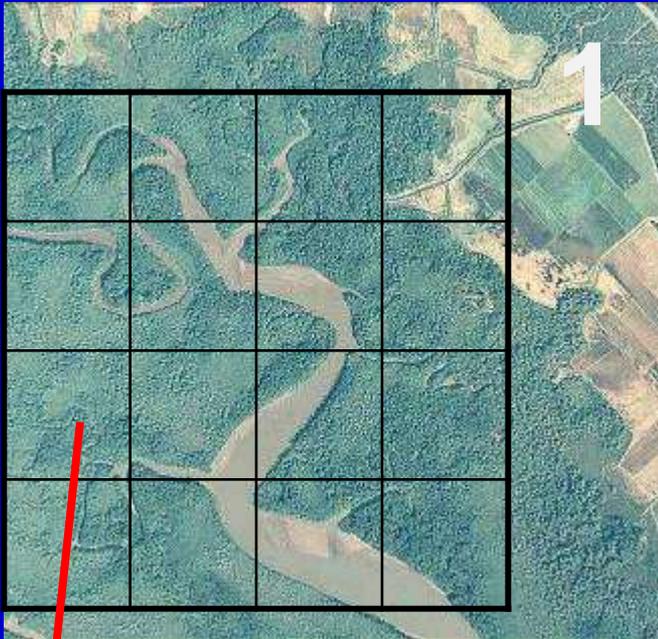


Foto TERRA: 23 dic. 97 Escala 1: 40.000

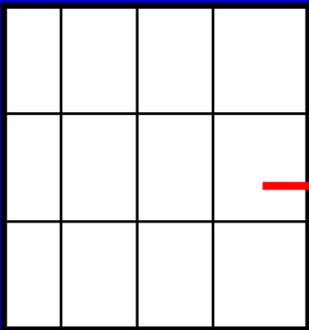
Muestreo Polietápico o Jerárquico

- Las unidades de muestreo de orden uno dan origen a las unidades de orden dos, éstas a su vez a las de orden tres y así sucesivamente hasta llegar al último elemento de la muestra. El diseño de muestreo es apropiado cuando se quiera estudiar una población a diferentes niveles de intensidad o detalle.
- **Ejemplo:** investigar producción de frutos de los aguacatillos. Este esquema de muestreo permite evaluar la variación a nivel nacional, regional, entre localidades en cada región y entre árboles en cada localidad
- El método es excelente para el estudio exhaustivo de poblaciones donde se requiera conocer los patrones de variabilidad natural, como por ejemplo en el área de mejoramiento genético y biosistemática

Ejemplo



- Nivel 1: Unidades de orden 1
- Nivel 2: Unidades de orden 2
- Nivel 3: Unidades de orden 3



Árboles, mantillo

Desventajas

- El tamaño de la muestra es usualmente pequeño en cada nivel o etapa.
- Las estimaciones no alcanzan generalmente la confiabilidad requerida (debido al reducido tamaño de la muestra).
- Se puede aumentar el tamaño de la muestra, aunque por la naturaleza misma del muestreo (jerárquico) esto significa un incremento considerable en los costos.

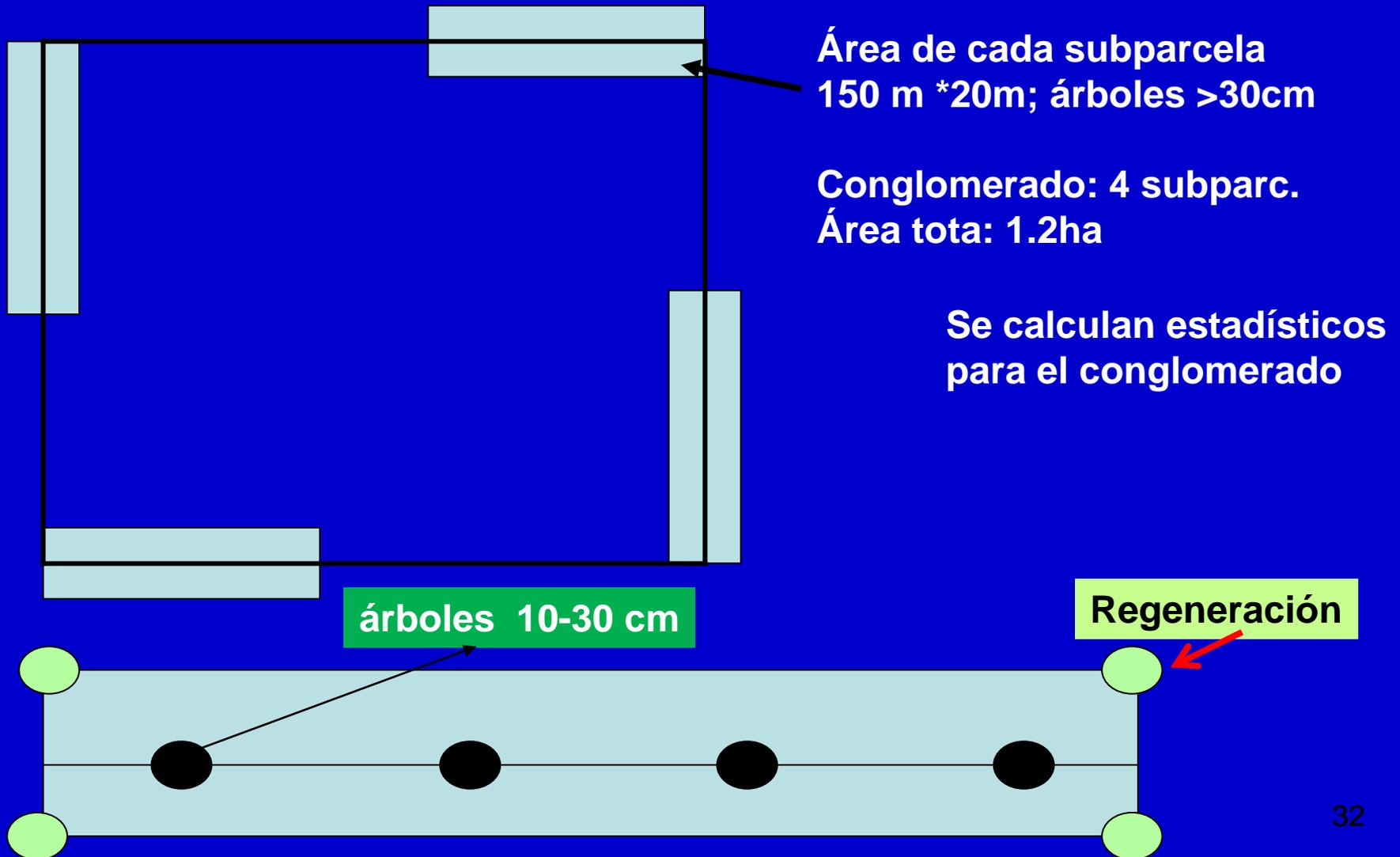
Muestreo Polifásico

- El muestreo polifásico utiliza algunas de las unidades en diferentes fases del muestreo y se utiliza cuando se requiere levantar datos a diferentes niveles de una población.
- Por ejemplo, el muestreo se denomina doble cuando se trabaja con dos fases y triple cuando se utilizan tres fases.
- Las fases pueden ser independientes o dependientes. Un ejemplo de muestreo triple es un inventario de cobertura que utiliza fotos aéreas en el levantamiento de vegetación (fase uno), parcelas en la fase de campo (fase dos) y árboles individuales para estimar producción de flores y frutos (fase tres).

Muestreo Conglomerados

- Unidades de muestreo se establecen en grupos o conjuntos.
- Cada grupo constituye un conglomerado.
- El método es especialmente útil para situaciones donde es difícil o muy costoso obtener una muestra independiente como por ejemplo en áreas boscosas de topografía escarpada y con escasas vías de comunicación.

Ejemplo: conglomerado



Muestreo sistemático

- Muestreo no probabilístico
- Requiere que el investigador (a) conozca la totalidad de los elementos que conforman la población.
- Seleccionar primera unidad de muestreo al azar y las remanentes siguiendo un patrón fijo y constante; por ejemplo, cada décima unidad.
- **Desventaja:** no permite calcular ni el error de muestreo ni el intervalo de confianza asociado a las estimaciones realizadas.

Ejemplo

	1	2	3	4	5	6	7	8	8	10
1	6,7	28,4	19,5	19,9	17,6	16,0	18,2	15,6	14,8	16,1
2	31,7	10,9	14,4	15,7	24,7	21,9	6,9	24,2	37,9	22,1
3	12,0	14,8	18,1	18,0	9,1	13,7	17,0	18,2	23,4	6,8
4	15,7	17,5	20,7	13,7	16,1	21,5	10,5	10,1	17,5	16,4
5	17,8	13,2	21,3	16,6	7,2	14,0	14,6	22,4	15,1	13,9
6	17,6	21,4	19,0	12,5	7,5	20,5	22,8	14,8	17,7	37,0
7	17,6	10,2	12,6	9,9	19,8	15,5	13,7	14,4	22,3	7,9
8	18,2	27,4	22,0	9,5	15,2	11,6	22,8	22,7	17,6	29,9
9	11,0	11,8	18,8	20,6	9,7	21,5	10,6	16,2	16,1	24,1
10	37,7	11,9	18,4	11,5	22,4	23,9	18,9	19,6	18,5	6,5

Población de observaciones

Selección de 1 en cinco con inicio aleatorio en la observación 2

Muestra: 28.14, 10.8. 14.8, 17.5, 13.2, 21.4.....



Un sitio adecuado para la practica de muestreo!!!!

Jorge Fallas
jfallas56@gmail.com