

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 standing on their hind legs and reaching out with their front paws towards each other, as if playing or interacting. The background is a blurred green field with some taller grass. The overall tone is warm and natural.

# AST0212 – 2016-1

Introducción al análisis de datos

Instituto de Astrofísica

Facultad de Física

Pontificia Universidad Católica de Chile

# Nuestro Semestre 2016-1

[illegible]

# Clase previa (Clase 9):

REPASO

1. Repaso de temas críticos de la clase previa ✓
  1. Correlación.
  2. Incerteza de parámetros en la correlación lineal. ✓
  3. Corrección de error sistemático. Extrapolación.
2. Coeficiente de correlación. ✓
3. Significación de diferencia en media ✓



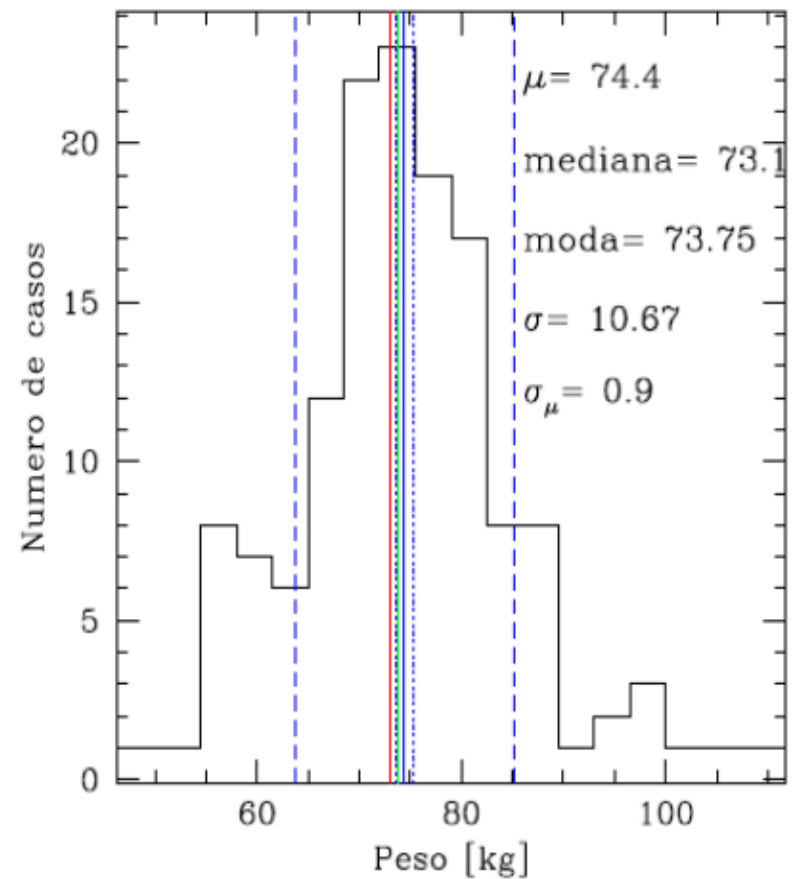
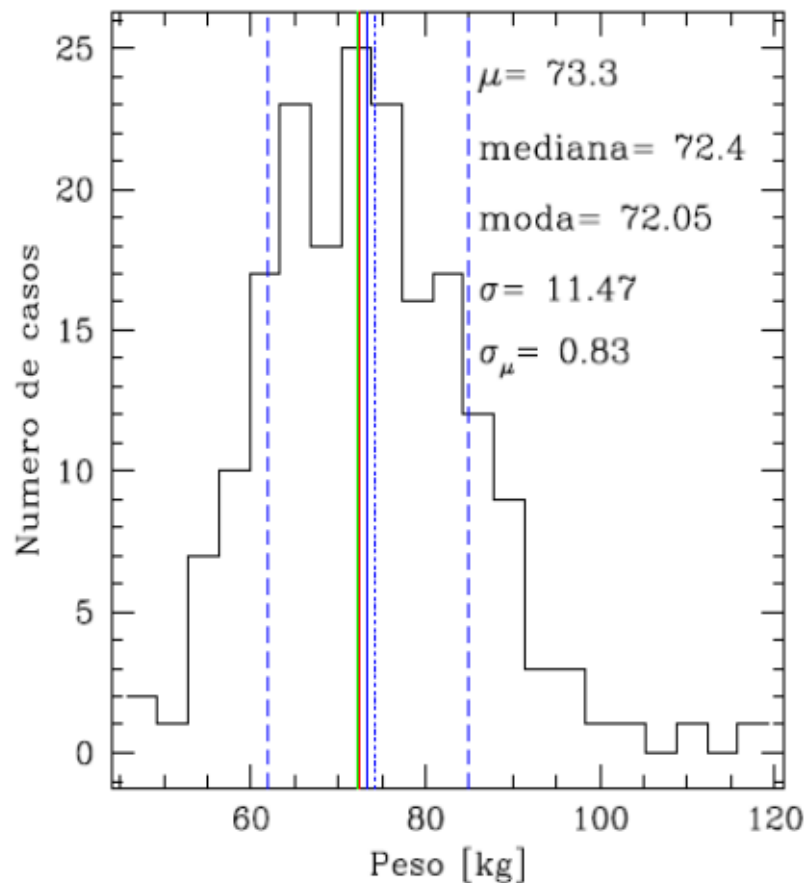
# Esta clase (Clase 10):

1. Repaso de temas críticos de la clase previa
  1. Significación de diferencia en media
  2. Correlación, coeficiente de correlación
2. Significación de un coeficiente de correlación.

# Comparación de dos distribuciones observadas: Significación de la diferencia de promedios

REPASO

G12345678-E-M-AC-peso.dat ;  $N_T = 190$  ; Bin = 3.5 kg | G12345678-E-M-DC-peso.dat ;  $N_T = 141$  ; Bin = 3.5 kg



# Comparación de dos distribuciones observadas:

## Significación de la diferencia de promedios

REPASO

Dados:

$$\bar{x}_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} x_{1,i}$$

$$\sigma_1 = \frac{1}{N_1 - 1} \sum_{i=1}^{N_1} x_{1,i}$$

$$\bar{x}_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} x_{2,i}$$

$$\sigma_2 = \frac{1}{N_2 - 1} \sum_{i=1}^{N_2} x_{2,i}$$

Tendremos el error del promedio:

¡Permite la estrategia de aumentar la muestra!

$$\sigma_{\bar{x}_1} = \frac{\sigma_1}{\sqrt{N_1}}$$

(Se obtienen de aplicar propagación de errores a las definiciones de  $\bar{x}_1$  y  $\bar{x}_2$ .)

$$\sigma_{\bar{x}_2} = \frac{\sigma_2}{\sqrt{N_2}}$$

Con estos elementos podemos construir el estimador  $t$ , con  $\nu = N_1 + N_2 - 2$  grados de libertad:

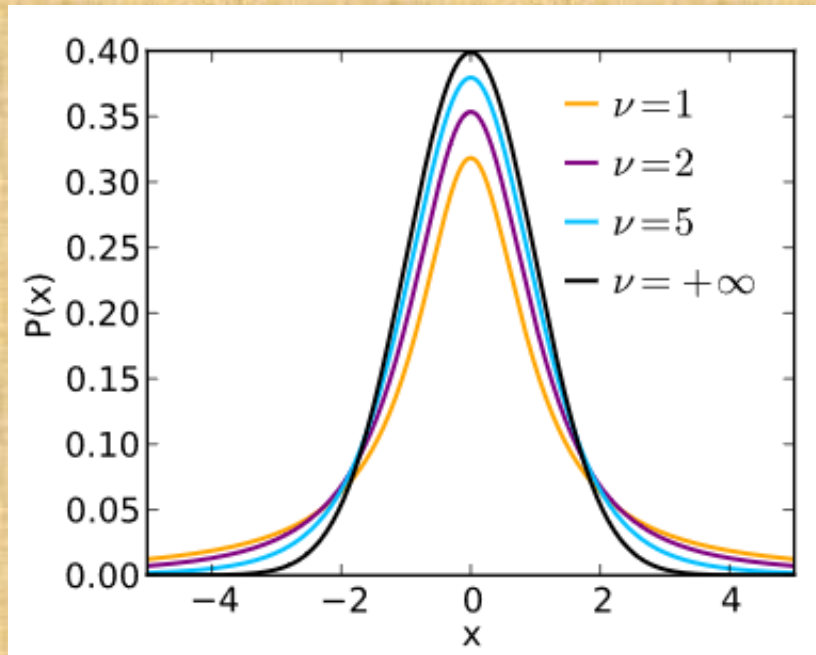
$$t = \frac{\bar{x}_1 - \bar{x}_2}{S_D}$$

donde

$$S_D = \sqrt{\frac{\sum_{N_1} (x_{1,i} - \bar{x}_1)^2 + \sum_{N_2} (x_{2,i} - \bar{x}_2)^2}{N_1 + N_2 - 2} \left( \frac{1}{N_1} + \frac{1}{N_2} \right)}$$

$S_D$  es el error estándar de la diferencia de promedios.  $t$  tiene FDP tipo  $t$  – Student.

# Distribución *t* de Student



La FDP de  $t$ ,  $A(t|\nu)$ , denota la probabilidad de que  $t$  sea, por azar, menor que el valor medido si los promedios  $\bar{x}_1$  y  $\bar{x}_2$  son realmente iguales. Un valor grande (por ejemplo 0.99) indica una *alta chance de medir un valor menor que el observado si  $\bar{x}_1 = \bar{x}_2$* . Esto es una indicación de que los promedios muy probablemente no sean los mismos. El valor complementario  $1 - A(t|\nu)$  es la probabilidad de medir un valor tan grande como  $t$  si  $\bar{x}_1 = \bar{x}_2$  (0.01 en el caso previo).

$$A(t|\nu) = \int_{-t}^t \frac{1}{\nu^{\frac{1}{2}} B\left(\frac{1}{2}, \frac{\nu}{2}\right)} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}} dx = \int_{-t}^t P(x) dx = 1 - I_{\frac{\nu}{\nu+t^2}}\left(\frac{\nu}{2}, \frac{1}{2}\right)$$

Donde  $B(a, b)$  es la función Beta, e  $I_x(a, b)$  es la función Beta incompleta (en este caso para  $x = \frac{\nu}{\nu+t^2}$ ,  $a = \nu/2$  y  $b = 1/2$ ).

Hay calculadores on-line para estas funciones, por ejemplo para  $1 - A(t|\nu)$ :

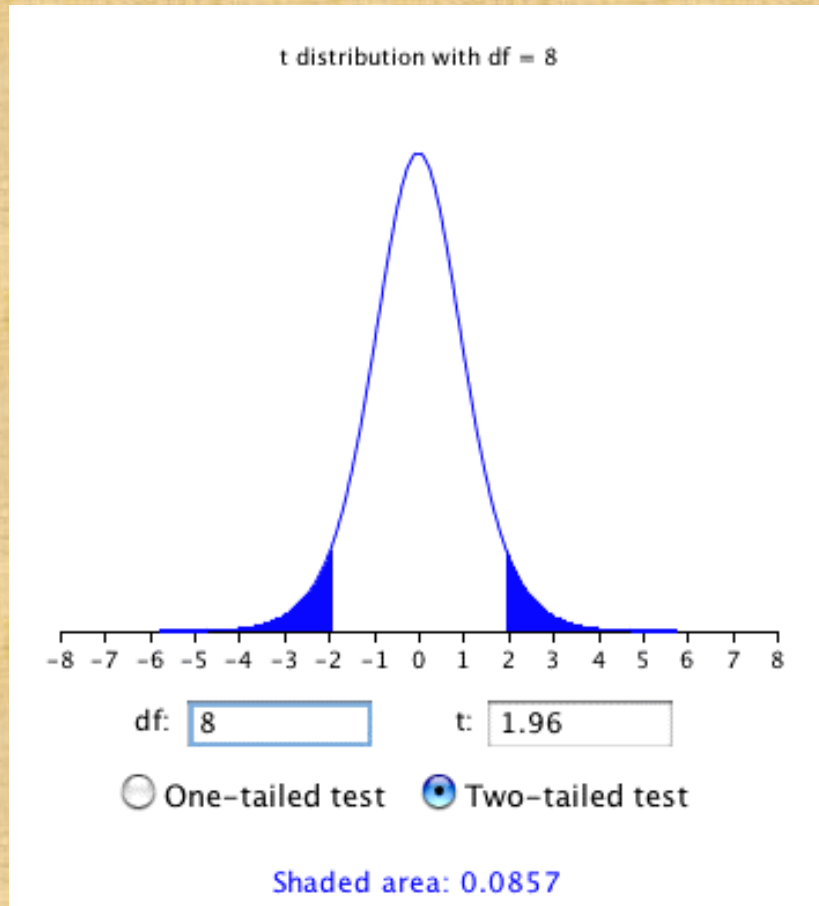
[http://onlinestatbook.com/2/calculators/t\\_dist.html](http://onlinestatbook.com/2/calculators/t_dist.html)



# Distribución *t* de Student

$$A(t|\nu) = \frac{1}{\nu^{\frac{1}{2}} B(\frac{1}{2}, \frac{\nu}{2})} \int_{-t}^t \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}} dx$$

La FDP *t* de Student es en realidad una FDP acumulativa (la integral entre  $-t$  y  $t$ ). Lo que estaba graficado en la lámina previa era el integrando.

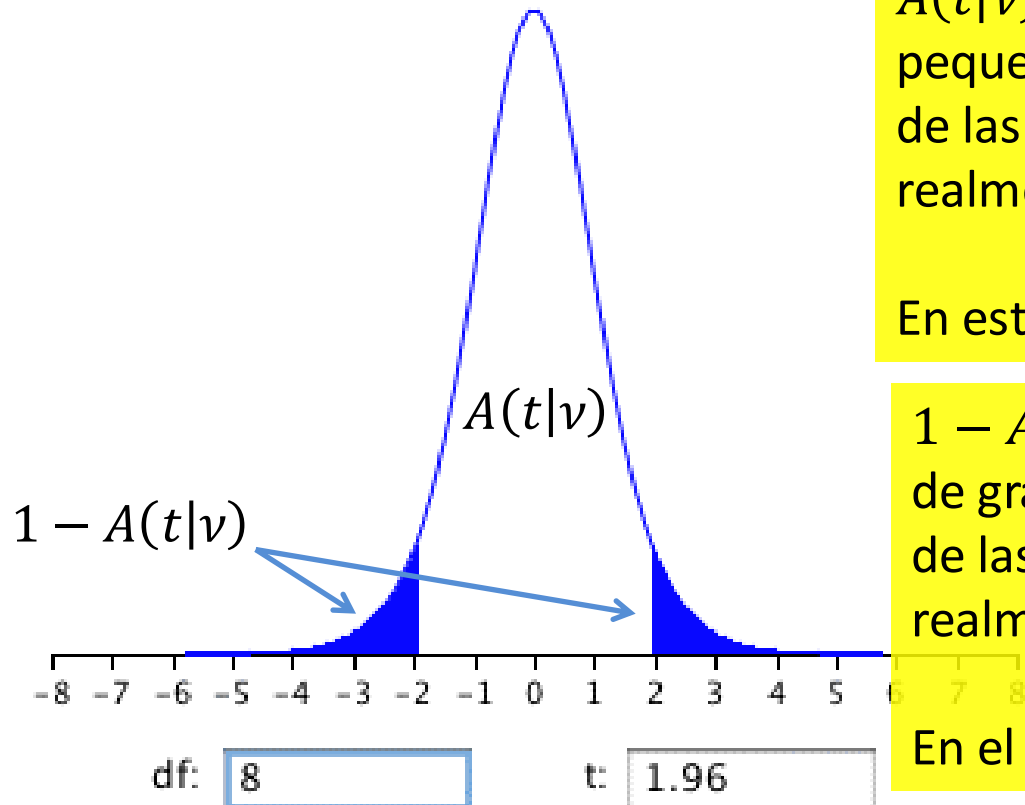


Este gráfico muestra  $A(t|\nu)$  como área blanca bajo la línea azul, y  $1 - A(t|\nu)$  como área azul en los extremos derecho e izquierdo de la distribución. Por simetría, debemos considerar ambas colas (ya que el orden en que hacemos la resta en la definición de  $t$  es arbitrario. La función  $1 - A(t|\nu)$  incorpora las dos colas naturalmente.



# Interpretación del *t de Student*

t distribution with df = 8



$A(t|v)$ : Probabilidad de que  $t$  sea así de pequeño por azar, cuando los promedios de las dos distribuciones comparadas son realmente iguales.

En este ejemplo:  $A(1.96|8) = 0.9143$

$1 - A(t|v)$ : Probabilidad de que  $t$  sea así de grande por azar, cuando los promedios de las dos distribuciones comparadas son realmente iguales.

En el ejemplo:  $1 - A(1.96|8) = 0.0857$

“Pequeño” o “grande” se miden en este contexto comparando con 1 ¿Se entiende?

# Coeficiente de correlación

REPASO

Partamos con un recordatorio de las ecuaciones de ajuste lineal de cuadrados mínimos:

$$y_i = ax_i + b$$

$$\sum_{i=1}^N \frac{y_i}{\sigma_i^2} = b \sum_{i=1}^N \frac{1}{\sigma_i^2} + a \sum_{i=1}^N \frac{x_i}{\sigma_i^2}$$

$$\sum_{i=1}^N \frac{x_i y_i}{\sigma_i^2} = b \sum_{i=1}^N \frac{x_i}{\sigma_i^2} + a \sum_{i=1}^N \frac{x_i^2}{\sigma_i^2}$$

$$\chi^2 = \sum_{i=1}^N \left( \frac{y_i - (ax_i + b)}{\sigma_i} \right)^2$$
$$\Delta = \sum_{i=1}^N \frac{1}{\sigma_i^2} \sum_{i=1}^N \frac{x_i^2}{\sigma_i^2} - \left( \sum_{i=1}^N \frac{x_i}{\sigma_i^2} \right)^2$$

$$a = \frac{1}{\Delta} \left( \sum_{i=1}^N \frac{1}{\sigma_i^2} \sum_{i=1}^N \frac{x_i y_i}{\sigma_i^2} - \sum_{i=1}^N \frac{x_i}{\sigma_i^2} \sum_{i=1}^N \frac{y_i}{\sigma_i^2} \right)$$

$$b = \frac{1}{\Delta} \left( \sum_{i=1}^N \frac{x_i^2}{\sigma_i^2} \sum_{i=1}^N \frac{y_i}{\sigma_i^2} - \sum_{i=1}^N \frac{x_i}{\sigma_i^2} \sum_{i=1}^N \frac{x_i y_i}{\sigma_i^2} \right)$$

Simplifiquemos para un caso sin  $\sigma$  ( es idéntico a imaginar  $\sigma = 1$  )

# Coeficiente de correlación (caso sin REPASO)

¿Tiene sentido la correlación  $y_i = ax_i + b$ ? Prestemos atención a la pendiente.

$$a = \frac{1}{\Delta} \left( N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i \right) = \frac{\Delta_s}{\Delta}$$
$$\Delta = N \sum_{i=1}^N x_i^2 - \left( \sum_{i=1}^N x_i \right)^2$$
$$x_i = a' y_i + b' \quad a' = \frac{\Delta_s}{\Delta'}$$
$$\Delta' = N \sum_{i=1}^N y_i^2 - \left( \sum_{i=1}^N y_i \right)^2$$

Si hay una correlación real entre  $x$  e  $y$  deberá existir una relación entre  $a, a', b$  y  $b'$ .

$$x_i = \frac{1}{a} y_i - \frac{b}{a} \Rightarrow a' = \frac{1}{a}; b' = -\frac{b}{a} \Rightarrow aa' = 1$$

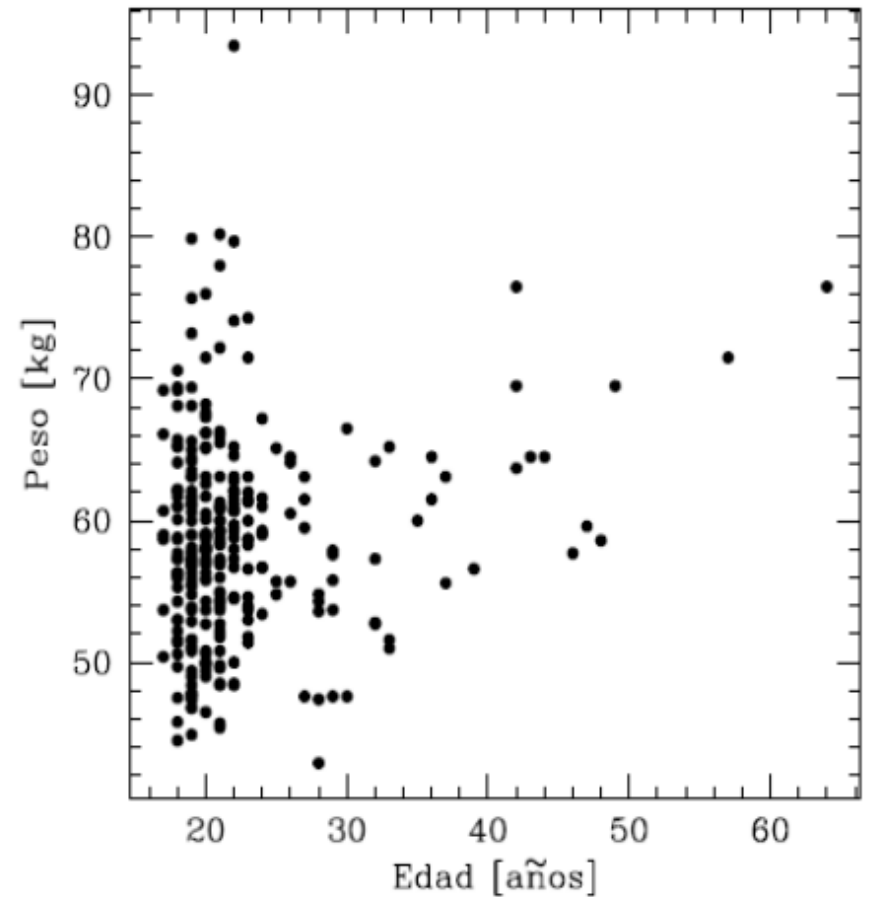
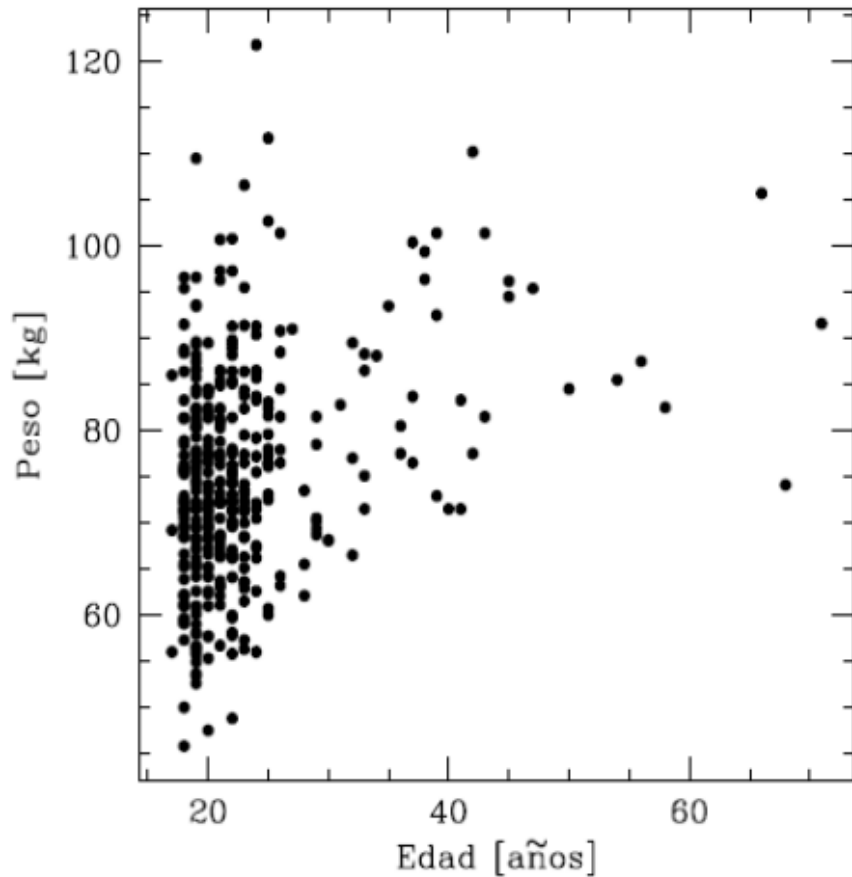
$$aa' = 1$$

Definimos  $r = \sqrt{aa'}$  cantidad llamada “coeficiente de correlación lineal”, que nos da una medida experimental del grado de correlación lineal, con valor entre 0 y  $\pm 1$ .

$$r = \frac{N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{\sqrt{N \sum_{i=1}^N x_i^2 - \left( \sum_{i=1}^N x_i \right)^2} \sqrt{N \sum_{i=1}^N y_i^2 - \left( \sum_{i=1}^N y_i \right)^2}}$$

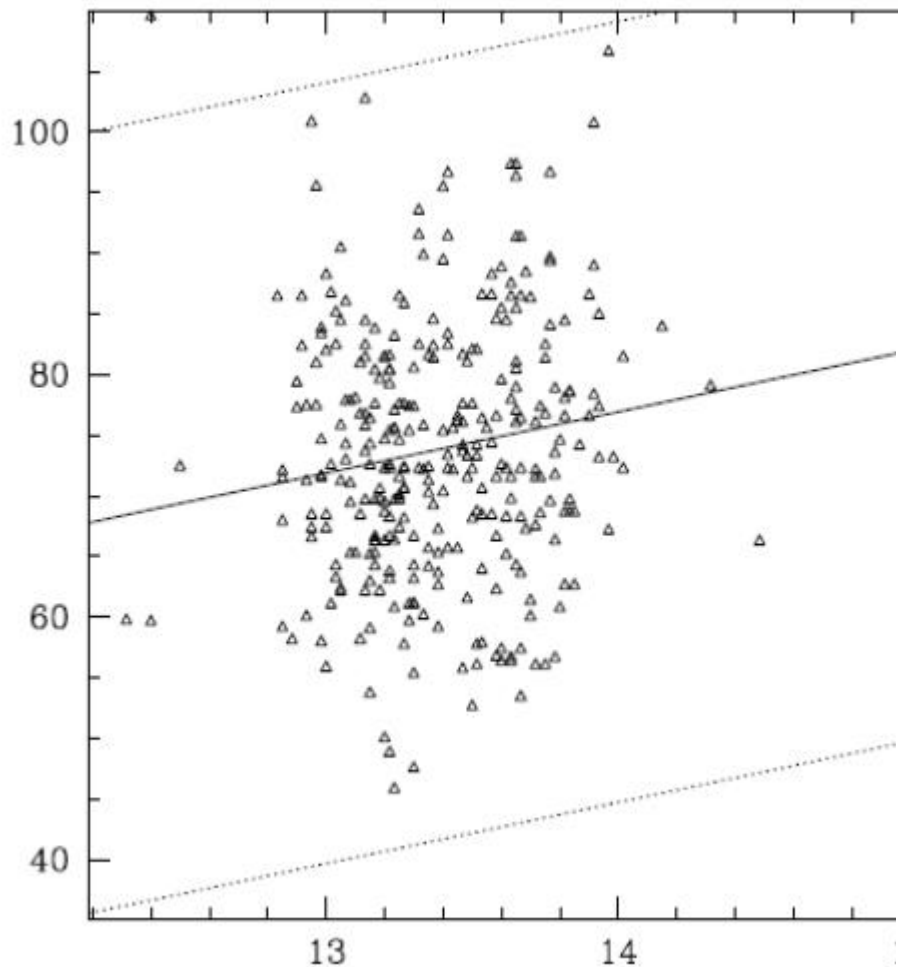
(Raimundo, antes de que tomáramos la raíz cuadrada, el numerador era un cuadrado)

# ¿Correlaciones casuales?





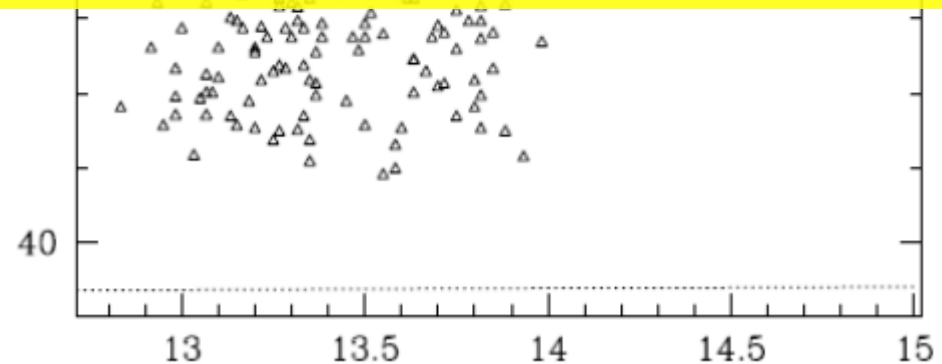
# ¿Correlaciones casuales?



En la hipótesis inicial (nula) de que  $x$  e  $y$  NO están correlacionados (y si los datos estás suficientemente agrupados, y si hay más que  $n \approx 20$  puntos) entonces  $r$  tiene distribución  $N(0, \sqrt{n})$ . En ese caso, la probabilidad de que  $r$  sea “así de grande” por azar (es decir, con  $x$  e  $y$  no correlacionadas). Está dada por:

$$1 - \operatorname{erf}\left(\frac{|r|\sqrt{n}}{\sqrt{2}}\right)$$

(esta  $\operatorname{erf}(x)$  es la real, no la del profe)



# Condoro del profe (2)

```
[aclocchi@localhost ~/source]$ more myFx.f
      double precision FUNCTION FF(X)
C
C adapted from erfcc in NR
C
C This one assumes that x is N(0,1) and returns F(x), the integral of a
C normalized gaussian distribution between -infty and x (i.e. c gives
C funcion P(x), eq. 26.2.2 in Abramowitz & Stegun). Note that
C
C          2.0d0*ff(x)-1.0
C
C gives the probability that the results fall between -x and x, as usual
C
C x = 1 =>  0.68268946714998124
C x = 2 =>  0.95449973904422536
C x = 3 =>  0.99730020379497075
C x = 4 =>  0.99993665751002569
C x = 5 =>  0.99999942669685327
C
C etc., etc., etc.
C
      implicit double precision (a-h,o-z)
C
      Z=DABS(X)/dsqrt(2.0d0)
      T=1./(1.+0.5*Z)
      ERFCC=T*DEXP(-Z*Z-1.26551223+T*(1.00002368+T*(.37409196+
*      T*(.09678418+T*(-.18628806+T*(.27886807+T*(-1.13520398+
*      T*(1.48851587+T*(-.82215223+T*.17087277))))))))))
      IF (X.LT.0.) ERFCC=2.-ERFCC
      FF=(2.0d0-ERFCC)/2.0d0
      RETURN
      END
```

Fin de ppt de Clase 10

# Preguntas:

En general lo están haciendo bastante mal con la tarea:

1. ¿Qué hace falta para que esto funcione mejor?
  - a) ¿Algo que pueda hacer desde la cátedra?
  - b) ¿Algo que podamos hacer desde las ayudantías?
  - c) ¿Algo que pueda hacer la universidad?