

# **Análisis y Visualización de Datos**

## **Clase 2**

## **Conceptos de Estadística Descriptiva**

**Docentes :**      **Soledad Palacios(UNLP)**

**Milagro Teruel (UNC)**

# Repaso de la clase anterior

- Estudiamos en profundidad el concepto de Probabilidad, estudiando sus tres variantes: probabilidad clásica, probabilidad frecuentista y probabilidad subjetiva
- Examinamos la clasificación de los datos: variables cuantitativas y variables categóricas
- Revisamos el concepto de variable aleatoria
- Analizamos las distribuciones de las variables aleatorias

# Definición de probabilidad

La **probabilidad** es una medida de la certidumbre asociada a un suceso o evento futuro y suele expresarse como un número entre 0 y 1 (o entre 0 % y 100 %).

Una forma tradicional de estimar algunas probabilidades sería obtener la frecuencia de un acontecimiento determinado mediante la realización de experimentos aleatorios, de los que se conocen todos los resultados posibles, bajo condiciones *suficientemente* estables.

# Estadística

La estadística es una rama de las matemáticas y una herramienta que estudia usos y análisis provenientes de una muestra representativa de datos, que busca explicar las correlaciones y dependencias de un fenómeno físico o natural, de ocurrencia en forma aleatoria o condicional.

El campo de la estadística tiene que ver con la recopilación, organización, análisis y uso de datos para tomar decisiones razonables basadas en tal análisis.

# Muestra y población

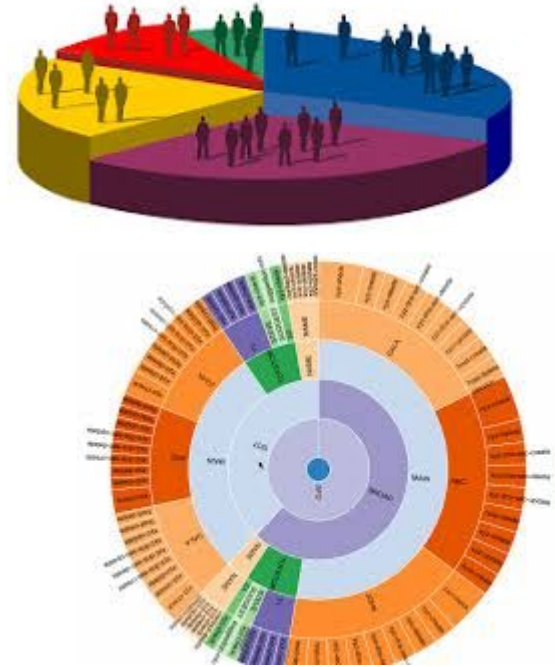
Cuando recogemos los datos muchas veces es imposible relevar la característica de interés de todos el grupo entero (**población**) o universo, se examina una pequeña parte del grupo, llamada muestra.



# Estadística Descriptiva

La parte de la estadística que estudia la muestra sin inferir alguna conclusión sobre la población es la estadística descriptiva.

En particular la estadística descriptiva trata sobre los métodos para recolectar, organizar y resumir datos.



# **Analicemos el dataset del Titanic**

**Vamos a trabajar con el dataSet del Titanic.**

**Veamos que datos tenemos**



# Miramos que datos tiene

les pedimos una descripción total a panda

In [5]:

```
titanic.describe()
```

Slide Type

Out[5]:

|       | PassengerId | Survived   | Pclass     | Age        | SibSp      | Parch      | Fare       |
|-------|-------------|------------|------------|------------|------------|------------|------------|
| count | 891.000000  | 891.000000 | 891.000000 | 714.000000 | 891.000000 | 891.000000 | 891.000000 |
| mean  | 446.000000  | 0.383838   | 2.308642   | 29.699118  | 0.523008   | 0.381594   | 32.204208  |
| std   | 257.353842  | 0.486592   | 0.836071   | 14.526497  | 1.102743   | 0.806057   | 49.693429  |
| min   | 1.000000    | 0.000000   | 1.000000   | 0.420000   | 0.000000   | 0.000000   | 0.000000   |
| 25%   | 223.500000  | 0.000000   | 2.000000   | 20.125000  | 0.000000   | 0.000000   | 7.910400   |
| 50%   | 446.000000  | 0.000000   | 3.000000   | 28.000000  | 0.000000   | 0.000000   | 14.454200  |
| 75%   | 668.500000  | 1.000000   | 3.000000   | 38.000000  | 1.000000   | 0.000000   | 31.000000  |
| max   | 891.000000  | 1.000000   | 3.000000   | 80.000000  | 8.000000   | 6.000000   | 512.329200 |

¿Qué hacemos con estos datos?



## Distribución de frecuencias para variables cuantitativas

Un tema íntimamente relacionado con los histogramas son las tablas de distribución de frecuencia, en definitiva los histogramas no son más que gráficos de tablas de distribución de frecuencia. La distribución de frecuencia de una variable cuantitativa consiste en un resumen de la ocurrencia de un dato dentro de una colección de categorías que no se superponen. Estas categorías las vamos a poder armar según nuestra conveniencia y lo que queramos analizar.

# Distribución de frecuencias

In [27]:

```
# Distribución de frecuencia.  
# 1ro creamos un rango para las categorías.  
contenedores = np.arange(0, 81., 10)  
  
# luego cortamos los datos en cada contenedor  
frec = pd.cut(titanic['Age'], contenedores)  
  
# por último hacemos el recuento de los contenedores  
# para armar la tabla de frecuencia.  
tabla_frec = pd.value_counts(frec)  
tabla_frec
```

Out[27]:

|              |     |
|--------------|-----|
| (20.0, 30.0] | 407 |
| (30.0, 40.0] | 155 |
| (10.0, 20.0] | 115 |
| (40.0, 50.0] | 86  |
| (0.0, 10.0]  | 64  |
| (50.0, 60.0] | 42  |
| (60.0, 70.0] | 17  |
| (70.0, 80.0] | 5   |

Name: Age, dtype: int64

Por defecto muestra el intervalo con más ocurrencias en primer lugar

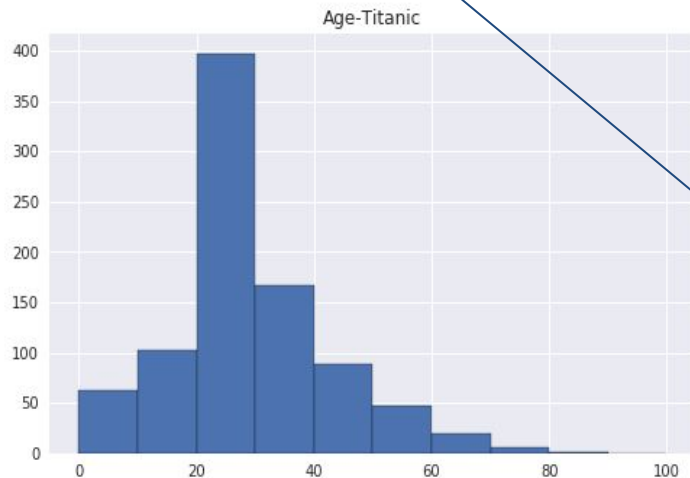
# Histograma

In [30]:

```
plt.title('Age-Titanic')
plt.hist(titanic.Age, edgeColor='black', bins=[0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100])
```

Slide Type

Out[30]: (array([ 62., 102., 397., 167., 89., 48., 19., 6., 1., 0.]),  
array([ 0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100])),  
<a list of 10 Patch objects>)



Con este parámetro seteamos el valor del intervalo a graficar

Con este parámetro seteamos la separación entre las barras

# Propiedades del histograma

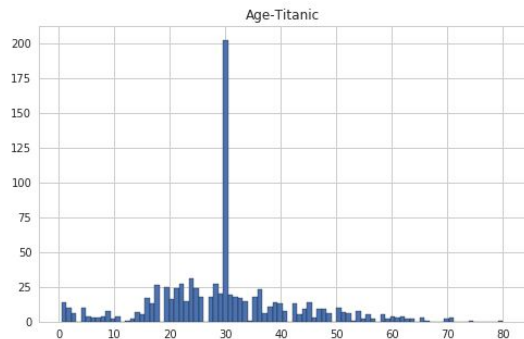
- Cada uno de los subintervalos representados por las barras son llamadas **clases**.
- En general se deben graficar entre 5 y 20 barras.
- El punto medio de cada clase es la marca de clase.
- La longitud de cada intervalo de clase es el ancho de clase
- el área de la barra debe ser proporcional a la frecuencia de la clase.

# Propiedades del histograma

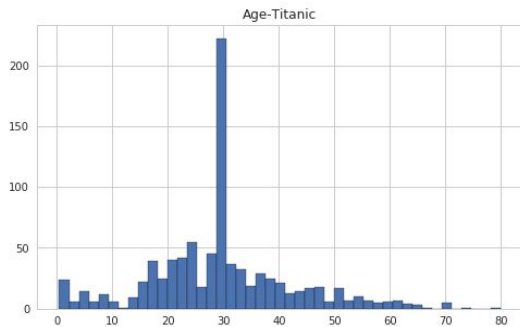
Los histogramas son útiles al proporcionar una impresión visual del aspecto que tiene la distribución de las mediciones, así como información sobre la dispersión de los datos. Al construir una tabla de frecuencias se pierde información, sin embargo esa pérdida de información es a menudo pequeña si se le compara con la facilidad de interpretación ganada al utilizar la distribución de frecuencias y el histograma.



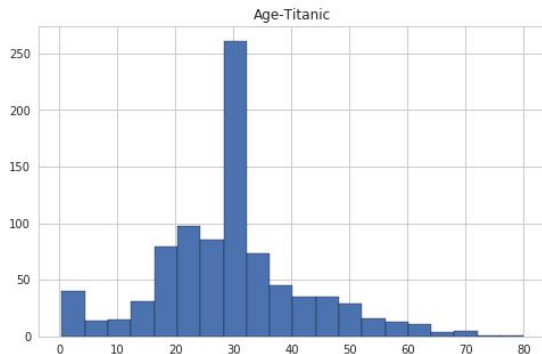
# Variación gráfica respecto las clases



```
plt.title('Age-Titanic')  
plt.hist(titanic.Age,  
edgeColor='black', bins=90)
```



```
plt.title('Age-Titanic')  
plt.hist(titanic.Age,  
edgeColor='black', bins=45)
```



```
plt.title('Age-Titanic')  
plt.hist(titanic.Age,  
edgeColor='black',  
bins=20)
```

# Gráficos de barras

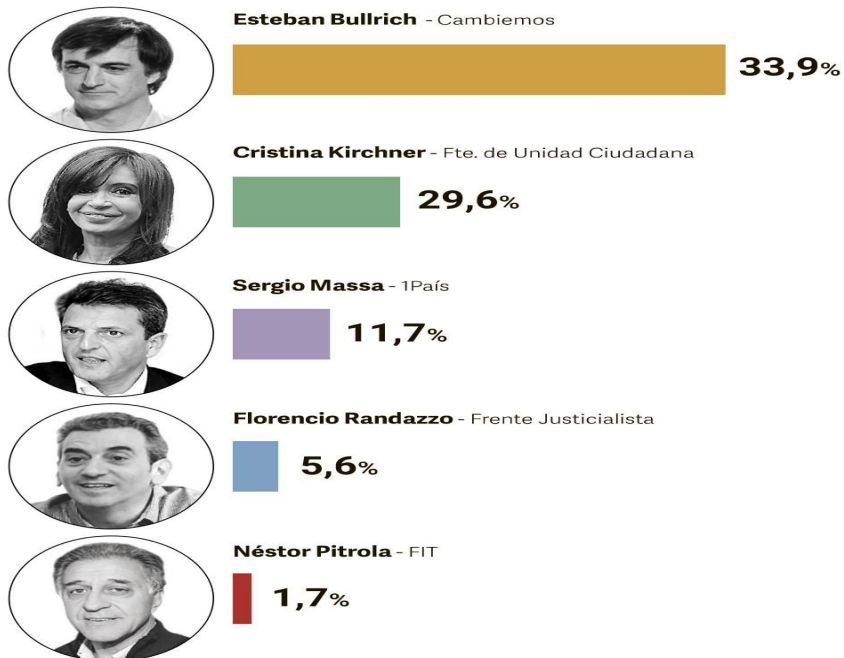
Las tablas de frecuencia también pueden emplearse en datos cualitativos o categóricos , es decir la muestra no consiste de valores numéricos (datos cuantitativos ) sino que los datos se ordenan en categorías y se registra cuántas observaciones caen en cada categoría (las categorías pueden ser masculino , femenino o fumador, no fumador o clasificar según nivel educativo: primario, secundario, terciario, universitario, ninguno). Cuando los datos son categóricos las clases se dibujan con el mismo ancho.



# Ejemplo

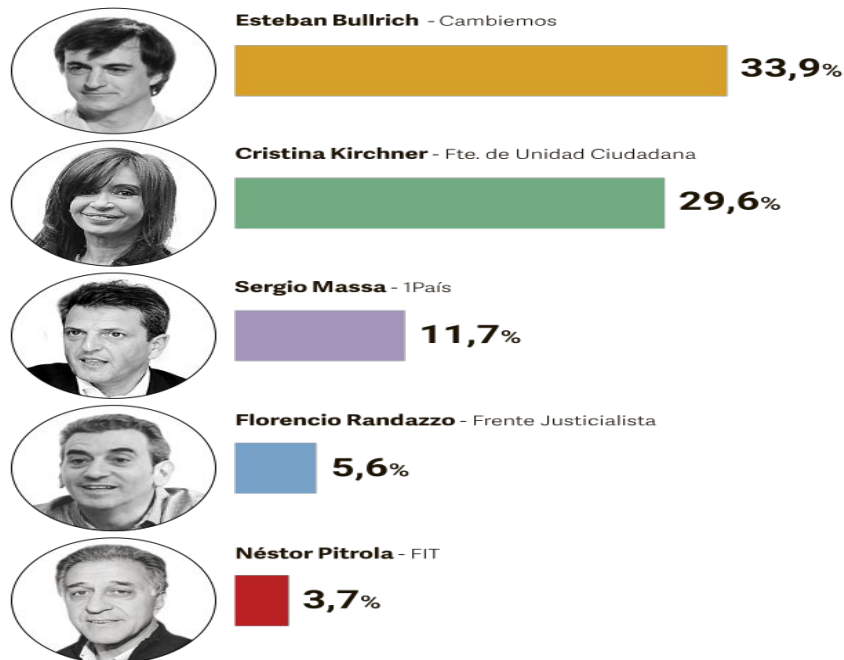
## En la Provincia de Buenos Aires

### INTENCIÓN DE VOTO A SENADOR NACIONAL



## En la Provincia de Buenos Aires

### INTENCIÓN DE VOTO A SENADOR NACIONAL





# Medidas Descriptivas

Del mismo modo que las gráficas pueden mejorar la presentación de los datos, las descripciones numéricas también tienen gran valor. Se presentan varias medidas numéricas importantes para describir las características de los datos.

Una característica importante de un conjunto de números es su localización o tendencia central .



# Medidas de localización - Media

La medida más común de localización o centro de un grupo de datos es el promedio aritmético ordinario o media. Ya que casi siempre se considera a los datos como una muestra, la media aritmética se conoce como media muestral.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

# Propiedades de la media

- Su cálculo es muy sencillo y en él intervienen todos los datos.
- Su valor es único para una serie de datos dada.
- Se usa con frecuencia para comparar poblaciones, aunque es más apropiado acompañarla de una medida de dispersión.

# Propiedades de la media

- Se interpreta como "punto de equilibrio" o "centro de masas"
- Minimiza las desviaciones cuadráticas de los datos respecto de cualquier valor prefijado, esto es, el valor de

$$\frac{1}{n} \sum_{i=1}^n (x_i - k)^2 \quad \text{es mínimo cuando} \quad k = \bar{x}$$

# Inconvenientes

- Para datos agrupados en intervalos (variables continuas) su valor oscila en función de la cantidad y amplitud de los intervalos que se consideren.
- Es una medida a cuyo significado afecta sobremanera la dispersión, de modo que cuanto menos homogéneos sean los datos, menos información proporciona
- En el cálculo de la media no todos los valores contribuyen de la misma manera. Los valores altos tienen más peso que los valores cercanos a cero

# Mediana

Representa el valor de la variable de posición central en un conjunto de datos **ordenados**.

Existen dos métodos para el cálculo de la mediana:

1. Considerando los datos en forma individual, sin agruparlos.
2. Utilizando los datos agrupados en intervalos de clase.



# Propiedades de la mediana

- Fácil de calcular si el número de observaciones no es muy grande.
- No se ve influenciada por valores extremos, ya que solo influyen los valores centrales.
- Se puede calcular para cualquier tipos de datos cuantitativos, incluso los datos con clase de extremo abierto.
- Es la medida de tendencia central más representativa en el caso de variables que solo admiten la escala ordinal.

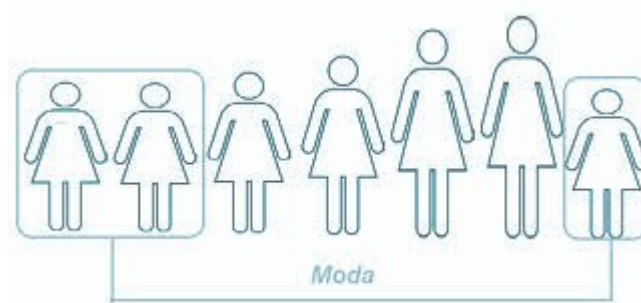
# Inconvenientes

- No utiliza en su “cálculo” toda la información disponible.
- No pondera cada valor por el número de veces que se ha repetido.
- Hay que ordenar los datos antes de determinarla.



# Moda

La moda es la observación que se presenta con mayor frecuencia en la muestra



# Propiedades de la moda

- No requiere cálculos.
- Puede usarse para datos tanto cuantitativos como cualitativos.
- Fácil de interpretar.
- No se ve influenciada por valores extremos.
- Se puede calcular en clases de extremo abierto.

# Inconvenientes

- Para conjuntos pequeños de datos su valor no tiene casi utilidad, si es que de hecho existe. Solo tiene significado en el caso de una gran cantidad de datos.
- No utiliza toda la información disponible.
- No siempre existe, si los datos no se repiten.
- En ocasiones, el azar hace que una sola observación no represente el valor más frecuente del conjunto de datos.
- Difícil de interpretar si los datos tiene 3 o más modas.

# El engañoso término medio

Martín Gardner en su libro “¡Aja! Paradojas” nos cuenta la historia de la fábrica del Sr. Artilugio (PRODILUGIO S.A.) y su nuevo empleado Félix, quien fuera víctima de los engaños de los parámetros estadísticos de posición. La dirección de PRODILUGIO está a cargo del Sr. Artilugio, su hermano y seis parientes. La fuerza laboral consiste en cinco encargados y diez operarios. Los negocios van bien y la fábrica precisa un operario más. El Sr. Artilugio entrevista a Félix, candidato al puesto, y le explica que su empresa paga muy bien, ya que el salario medio es de \$600 semanales. Al cabo de unos cuantos días, Félix quiso ver al jefe. Este fue su diálogo: Félix: -¡Me ha engañado usted! He hablado con los otros operarios y ninguno gana más de \$200 a la semana. ¿Cómo puede ser que el salario medio sea de \$600? Sr. Artilugio: -Vamos Félix, no se excite. El salario medio es de \$600 y se lo voy a demostrar. Vea esta nómina por favor.

# El engañoso término medio

| Empleado                          | Sueldo semanal  |
|-----------------------------------|-----------------|
| Sr. Artilugio                     | \$4800          |
| Hermano del Sr. Artilugio         | \$2000          |
| 6 parientes                       | \$500 (c/u)     |
| 5 capataces                       | \$400 (c/u)     |
| 10 operarios                      | \$200 (c/u)     |
| <b>Total</b> <b>23 empleados:</b> | <b>\$13.800</b> |

# El engañoso término medio

El sueldo promedio **resulta entonces:**  $\frac{\$13800}{23} = \$600$

Félix consideró que, si bien era cierto que el sueldo medio resultaba de \$600, de todas formas lo habían engañado.

Pero el Sr. Artilugio insistió y le propuso ordenar los sueldos de mayor a menor:

4800, 2000, 500, 500, 500, 500, 500, 500, 400, 400, 400, **400**, 400, 200, 200, 200, 200, 200, 200, 200, 200.

# El engañoso término medio

El valor 400 indicado en rojo, ocupa **el lugar central** en esa sucesión (existe la misma cantidad de valores de la variable sueldo por encima que por debajo de él). El Sr. Artilugio le explicó a Félix que esa era la **mediana**.

Pero Félix continuaba insatisfecho con estos argumentos y volvió a preguntar ya un poco alterado:

Félix: -¿Y qué significan entonces los \$200?

Sr. Artilugio: -Eso, muchacho, se llama **moda**. Y es el salario ganado por el **mayor número de personas**.

Sin duda, después de la experiencia de Félix, la representatividad de la media, la moda y la mediana para caracterizar una situación (en este caso, la situación salarial de los empleados del Sr. Artilugio) será un motivo de análisis y discusión.

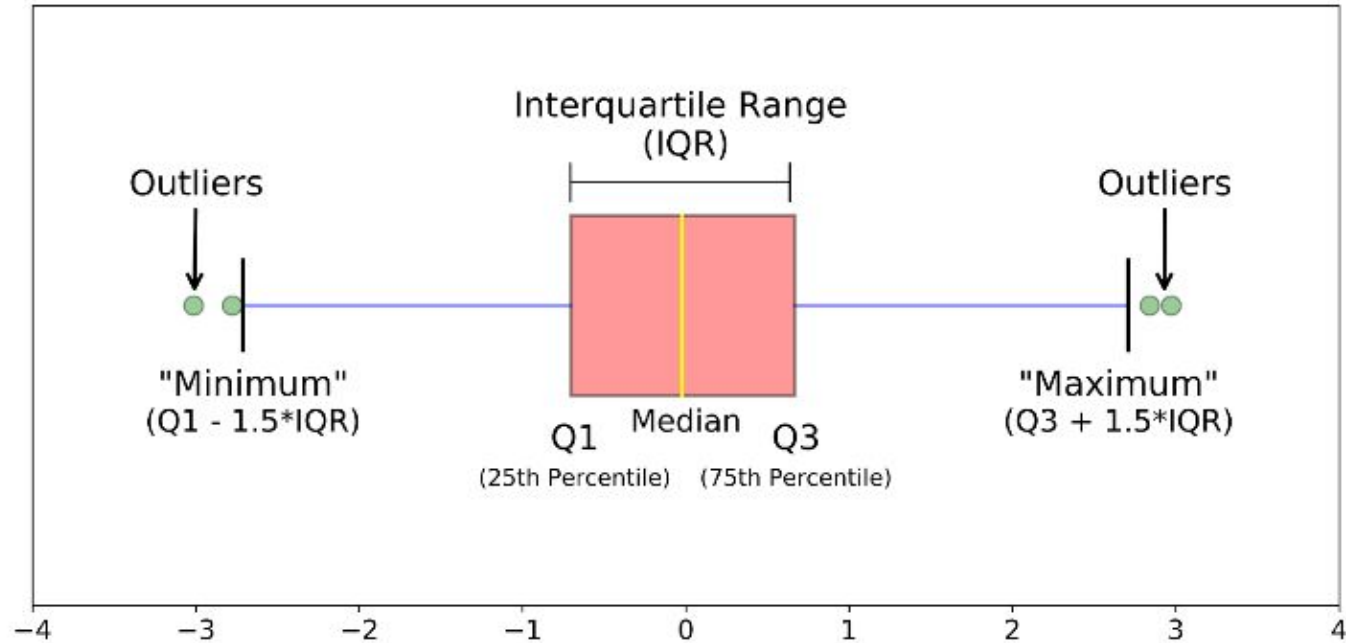
# percentiles

El **percentil** es una medida de posición usada en estadística que indica que una vez ordenados los datos de menor a mayor, el valor de la variable por debajo del cual se encuentra un porcentaje dado de observaciones en un grupo de observaciones.

[tablas de crecimiento de niños](#)



# boxplots



# Medidas de variabilidad

La localización o tendencia central no necesariamente proporciona información suficiente para describir datos de manera adecuada.

Las medidas de variabilidad nos informan sobre el grado de concentración o dispersión que presentan los datos respecto a su promedio.



# Clasificación

## Medidas dimensionales

- Rango
- Rango intercuartílico
- Varianza
- Desviación típica
- covarianza

## Medidas adimensionales

- Coeficiente de Pearson
- Rango intercuartílico
- Varianza
- Desviación típica
- Covarianza

# Rango de la muestra y rango intercuartílico

Una medida muy sencilla de variabilidad es el rango de la muestra , definido como la diferencia entre las observaciones más grande y más pequeña. Es decir

$$R = \text{Max } x_i - \text{Min } x_i$$

# Inconvenientes del rango

- El rango ignora toda la información que hay en la muestra entre las observaciones más chica y más grande.
- Además es muy sensible a los datos de los extremos por lo cual en algunas ocasiones es recomendable usar el rango intercuartílico.



# Desviación estándar y Varianza Muestral

Un modo de medir la variabilidad de los datos de una muestra sería tomar algún valor central, por ejemplo la media, y calcular el promedio de las distancias a ella. Mientras mayor sea este promedio, más dispersión deberían presentar los datos.

Sin embargo, esta idea no resulta útil, ya que las observaciones que se encuentran a la derecha de la media tendrán distancias (o desviaciones) positivas, en tanto que las observaciones menores que la media tendrán distancias negativas y la suma de las distancias a la media será inevitablemente igual a cero. Un modo de evitar este inconveniente es elevar las distancias al cuadrado y de este modo tener todos sumandos positivos.

# Desviación estándar y Varianza Muestral

La desviación estándar mide cuán lejos se encuentran los datos de la media muestral.

Definimos la varianza de una muestra de observaciones  $X_1$ ,  $X_2$ , ...,  $X_n$ , cuya media es  $\bar{X}$ , como

$$S^2 = \frac{\sum_i (X_i - \bar{X})^2}{n}$$

Porqué elevamos al cuadrado el numerador?

# Desviación estándar y Varianza Muestral

También se define con la siguiente

Corrección de Bessel

$$S^2 = \frac{\sum_{j=1}^n (X_j - \bar{X})^2}{n-1}$$

Definida con  $(n - 1)$  en el denominador la varianza muestral posee una propiedad deseable, resulta ser insesgado, esto es, en promedio no subestima ni sobrestima el valor de la varianza poblacional. Revisaremos esto en la 3° clase.



# Desviación estándar

La varianza muestral puede pensarse como “promedio” de las distancias a la media al cuadrado.

Sin embargo, la varianza no tiene las mismas unidades que los datos. Para salvar este inconveniente, definimos la desviación estándar muestral como la raíz cuadrada positiva de la varianza

$$s = \sqrt{s^2} = \sqrt{\frac{\sum (X - \bar{X})^2}{n - 1}}$$

# Volviendo al problema de los sueldos

Recordemos como se distribuían los sueldos en la empresa del Sr. Artilugio:

| Empleado                  | Sueldo semanal  |
|---------------------------|-----------------|
| Sr. Artilugio             | \$4800          |
| Hermano del Sr. Artilugio | \$2000          |
| 6 parientes               | \$500 (c/u)     |
| 5 capataces               | \$400 (c/u)     |
| 10 operarios              | \$200 (c/u)     |
| <b>Total</b>              | <b>\$13.800</b> |

# Volviendo al problema de los sueldos

Allí, el salario promedio era de \$600, sin embargo, Félix (y nosotros) vimos que no era representativo dada la gran dispersión de valores de sueldo que existe con respecto a esa media. Ahora vamos a calcular esa dispersión

$$\sigma^2 = \frac{(4800-600)^2 + (2000-600)^2 + (500-600)^2 \cdot 6 + (400-600)^2 \cdot 5 + (200-600)^2 \cdot 10}{23}$$

$$\sigma^2 \cong 933.043,5 \Rightarrow \sigma = \sqrt{933.043,5} \Rightarrow \sigma = 965.94$$

Este valor tan grande de la desviación estándar (incluso superior al promedio) pone aun más en evidencia la falta de representatividad del sueldo medio de \$600 en la empresa del Sr. Artilugio. Cualquier distribución de sueldos que usted proponga con menor dispersión de valores con respecto a la media (esto es con un desvío estándar menor) tendrá una media más representativa de la situación que la de esta empresa.

# tendencias

La tendencia habitual si se tiene una variable descrita en los términos de la Media $\pm$ Desviación estándar es a hacer aquellas típicas inferencias que sólo son ciertas si la variable se ajusta bien a la distribución normal:

- $M \pm 1DE$  supone el 68.5% aproximadamente de la población,
- $M \pm 2DE$  supone el 95% aproximadamente de la población
- $M \pm 3DE$  supone el 99.5% aproximadamente de la población

# tendencias

Si la variable no se ajusta a una distribución normal esas inferencias en absoluto son ciertas. Para evitar esta muy habitual inferencia inconsciente es mejor trabajar, evidentemente, en estos casos de no ajuste a la normalidad, con la Mediana y el Rango intercuartílico que son medidas que digamos están más próximas a la descripción propiamente dicha y no tienen tantas connotaciones inferenciales como las tienen la Media y la Desviación estándar.



# tendencias

El uso de unos u otros descriptores no depende del tamaño muestral, depende de la normalidad de la muestra, de su ajuste a la campana de Gauss.



# Asimetría o skewness

La asimetría es la medida que indica la simetría de la distribución de una variable respecto a la media aritmética, sin necesidad de hacer la representación gráfica. Los coeficientes de asimetría indican si hay el mismo número de elementos a izquierda y derecha de la media.



# Asimetría o skewness

- **Asimetría negativa:** la cola de la distribución se alarga para valores inferiores a la media.
- **Simétrica:** hay el mismo número de elementos a izquierda y derecha de la media. En este caso, coinciden la media, la mediana y la moda. La distribución se adapta a la forma de la campana de Gauss, o distribución normal.
- **Asimetría positiva:** la cola de la distribución se alarga para valores superiores a la media.



# Kurtosis y asimetría

También medidas que indican de la simetría o asimetría de la distribución y del achatamiento o no de la misma.

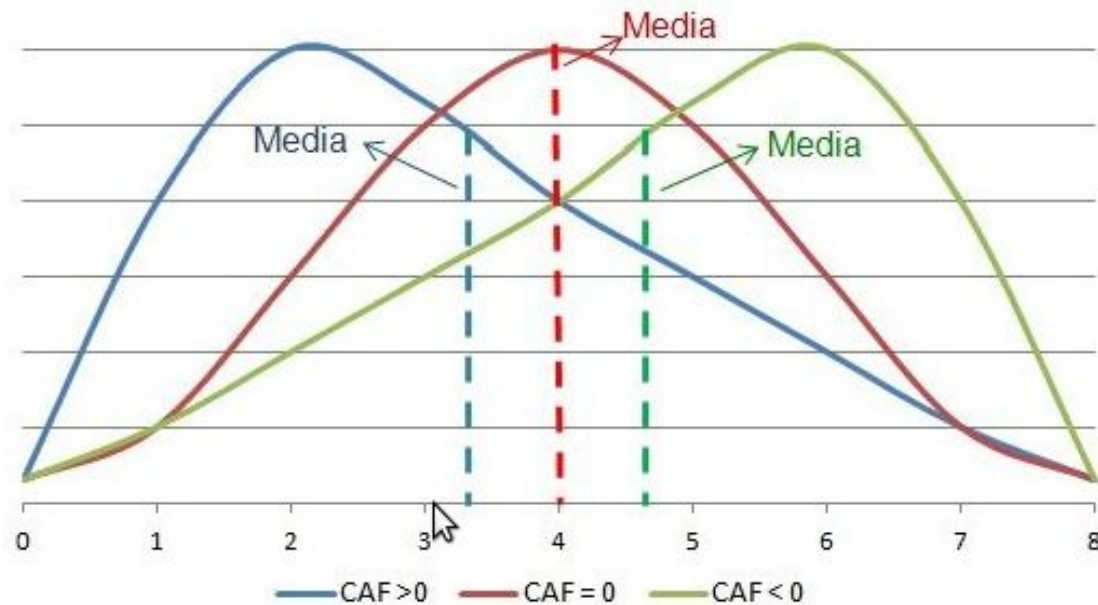
Empezando con la simetría, es lógico pensar que si la distribución tiene una única moda y es simétrica, entonces las tres medidas de centralización coinciden. Si no es simétrica, suele suceder que la mediana esté comprendida entre la moda y la media.

# Asimetría

El **coeficiente de asimetría de Fisher**  $CA_F$  evalúa la proximidad de los datos a su media  $\bar{x}$ . Cuanto mayor sea la suma  $\sum (x_i - \bar{x})^3$ , mayor será la asimetría. Sea el conjunto  $X = (x_1, x_2, \dots, x_N)$ , entonces la fórmula de la asimetría de Fisher es:

$$CA_F = \frac{\sum_{i=1}^N (x_i - \bar{x})^3}{N \cdot S_x^3}$$

# Asimetría



# Curtosis o apuntamiento

La **curtosis** de una variable estadística/aleatoria es una característica de forma de su distribución de frecuencias/probabilidad.

Una curtosis grande implica una mayor concentración de valores de la variable tanto muy cerca de la media de la distribución (pico) como muy lejos de ella (colas), al tiempo que existe una relativamente menor frecuencia de valores intermedios. Esto explica una forma de la distribución de frecuencias/probabilidad con colas más gruesas, con un centro más apuntado y una menor proporción de valores intermedios entre el pico y colas.

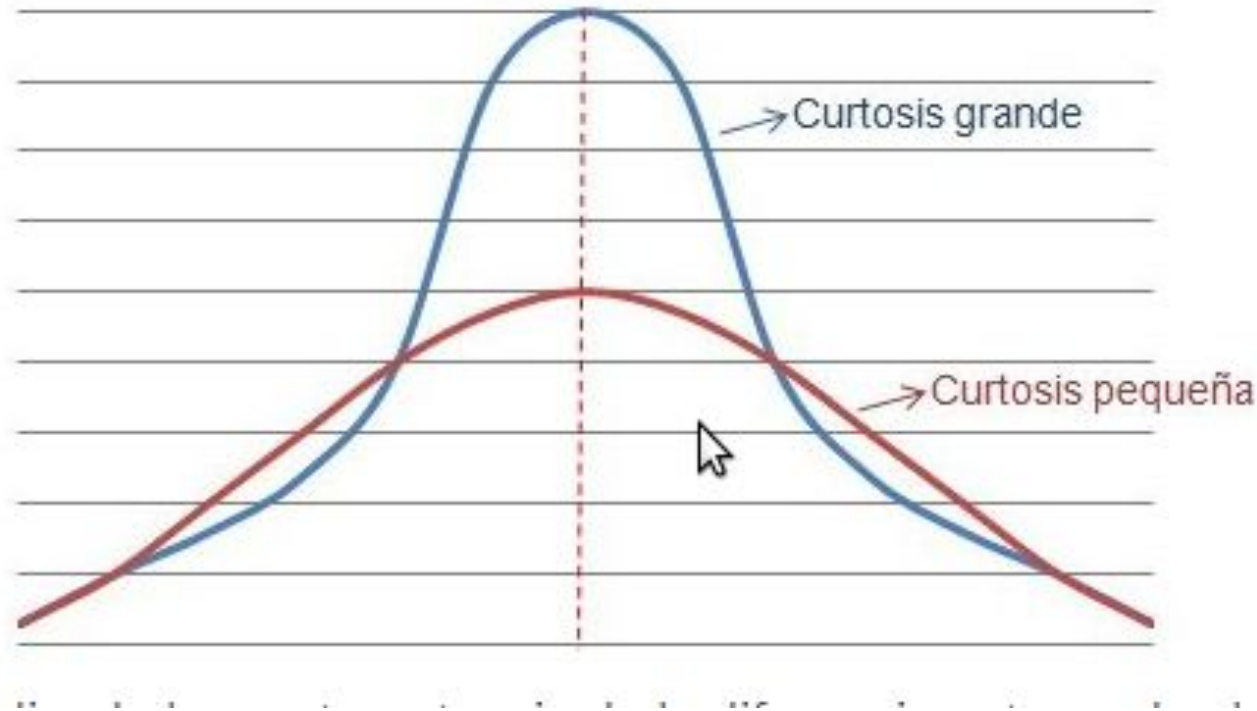
# Curtosis

Un coeficiente de apuntamiento o de curtosis es el cuarto momento con respecto a la media estandarizado que se define como:

$$Curtosis = \frac{\sum_{i=1}^N (x_i - \bar{x})^4}{N \cdot S_x^4} - 3$$

siendo  $\bar{x}$  la media y  $S_x$  la desviación típica

# Curtosis



**Una mayor  
curtosis no  
implica una  
mayor  
varianza, ni  
viceversa.**

# Pruebas paramétricas

## Supuestos de las paramétricas

- Normalidad
- Homocedasticidad: Las varianzas de los distintos grupos deben ser iguales, hay homogeneidad de varianzas.
- Respecto a los errores:
  - Los errores son independientes entre si
  - Se distribuyen normalmente dentro de cada población del grupo

# Pruebas no paramétricas

- Métodos de distribución libres, no requiere conocer la distribución de la muestra
- Permiten estudiar la forma de la distribución de la población de la que se extrajo la muestra
- Se denominan pruebas de bondad de ajuste y se usan para contrastar si los datos de la muestra proceden de cierta distribución.



# Bondad de ajuste

Resumen la discrepancia entre los valores observados y los valores esperados en el modelo de estudio.

Dentro de los test mas usados:

- Test de Kolmogorov-Smirnov (Test KS)
- Gráficos QQ

# Test de KS

## Ejemplo

Se quiere probar la hipótesis que una determinada distribución es normal.

Los valores ordenados para una muestra de tamaño 10 para esta distribución son:

66,72,81,94,112,116,124,140,145,155

¿qué conclusión puede obtenerse?



# Test de KS

¿Que nos vamos a preguntar?

$H_0$ : Los datos analizados siguen una distribución M.

$H_1$ : Los datos analizados no siguen una distribución M.

Analíticamente el proceso utiliza la siguiente transformación

$z = (\text{valor observado} - \text{media}) / \text{varianza}$

# Test de KS

Es entonces que el algoritmo a seguir para calcular lo que denominamos estadístico es el siguiente:

1. Ordenar las observaciones de mayor a menor
2. Acumulamos las equiprobabilidades  $1/N$  siendo  $N$  el orden que le haya tocado a cada observación
3. En una tercera columna anotamos la puntuación  $z$  de cada observación

# Test de KS

Armamos la tabla

| $j$           | valores | $F(j/n)$ | $\frac{j}{n} - F\left(\frac{j}{n}\right)$ | $\frac{j-1}{n} - F\left(\frac{j}{n}\right)$ |
|---------------|---------|----------|---|---|
| 1             | 66      | 0,48     | -0,38                                     | 0,48  |
| 2             | 72      | 0,51     | -0,31                                     | 0,41  |
| 3             | 81      | 0,56     | -0,26                                     | 0,36  |
| 4             | 94      | 0,61     | -0,21                                     | 0,31  |
| 5             | 112     | 0,67     | -0,17                                     | 0,27  |
| 6             | 116     | 0,69     | -0,09                                     | 0,19  |
| 7             | 124     | 0,71     | -0,01                                     | 0,11  |
| 8             | 140     | 0,75     | 0,05                                      | 0,05  |
| 9             | 145     | 0,77     | 0,13                                      | -0,03                                       |
| 10            | 155     | 0,79     | 0,21                                      | -0,11                                       |
| $d = 0,48315$ |         |          |   |   |

Encontramos el máximo y si el máximo es mayor que el nivel de significancia que queremos entonces no podemos rechazar  $H_0$

# Test de KS

```
In [30]: 1 from scipy import stats
          2 import numpy as np
          3
          4
          5 x = np.array([66, 72, 81, 94, 112, 116, 124, 140, 145, 155])
          6
          7 stats.kstest(x, 'norm')
          8
          9
         10
```

```
Out[30]: KstestResult(statistic=1.0, pvalue=0.0)
```

Al dar el p-value 0 no podemos decir que la distribución pertenezca a una distro normal



# Correlación

En muchas ocasiones es necesario estudiar conjuntamente dos características de un fenómeno aleatorio, es decir, el comportamiento conjunto de dos variables aleatorias, intentando explicar la posible relación existente entre ellas.

La correlación trata de establecer la relación o dependencia que existe entre las dos variables que intervienen en una distribución bidimensional.



# Relaciones entre variables numéricas

La existencia de algún tipo de asociación entre dos o más variables representa la presencia de algún tipo de tendencia o patrón de emparejamiento entre los distintos valores de esas variables.

Complementariamente, se habla de independencia entre variables cuando no existe tal patrón de relación entre los valores de las mismas.





# Tablas de Contingencia

Cuando se trabaja con variables categóricas, los datos suelen organizarse en tablas de doble entrada en las que cada entrada representa un criterio de clasificación (una variable categórica). Como resultado de esta clasificación, las frecuencias (el número o porcentaje de casos) aparecen organizadas en casillas que contienen información sobre la relación existente entre ambos criterios. A estas tablas de frecuencias se les llama **tablas de contingencia**.

# Tablas de Contingencia

|               | <b>Diestro</b> | <b>Zurdo</b> | <b>TOTAL</b> |
|---------------|----------------|--------------|--------------|
| <b>Hombre</b> | 43             | 9            | 52           |
| <b>Mujer</b>  | 44             | 4            | 48           |
| <b>TOTAL</b>  | 87             | 13           | 100          |

# Tablas de Contingencia

```
In [9]: # Tabla de contingencia class / survived  
pd.crosstab(index=titanic['survived'],  
             columns=titanic['class'], margins=True)
```

Out[9]:

| class    | 1st class | 2nd class | 3rd class | All  |
|----------|-----------|-----------|-----------|------|
| survived |           |           |           |      |
| no       | 122       | 167       | 528       | 817  |
| yes      | 203       | 118       | 178       | 499  |
| All      | 325       | 285       | 706       | 1316 |

# Relación entre variables

En las técnicas de relación el objetivo básico es detectar relación entre variables. La focalización está puesta, en este tipo de técnicas, en las variables, no en las poblaciones que pueda haber en el estudio. Los protagonistas son las variables y detectar covariación entre ellas, detectar que la variación que vemos en una de ellas tiene conexión con la de la otra.

# Variables aleatorias independientes

Intuitivamente decimos que dos variables,  $X$  e  $Y$ , son independientes si el valor que toma una de ellas no influye de ninguna manera sobre el valor que toma la otra. Esto lo establecemos más formalmente:

Sea  $(X, Y)$  una variable aleatoria bidimensional discreta. Sea  $p(x_i, y_j)$  su fdp conjunta y  $p(x_i)$  y  $q(y_j)$  las correspondientes fdp marginales de  $X$  e  $Y$ . Decimos que  $X$  e  $Y$  son variables aleatorias independientes si y sólo si

$$p(x_i, y_j) = p(x_i) q(y_j) \quad \forall (x_i, y_j) \in R_{XY}$$

# Covarianza

Es un valor que indica el grado de variación conjunta de dos variables aleatorias respecto a sus medias. Es el dato básico para determinar si existe una dependencia entre ambas variables y además es el dato necesario para estimar otros parámetros básicos, como el coeficiente de correlación lineal o la recta de regresión.

$$Cov(X, Y) = \frac{\sum_1^n (x_i - \bar{x})(y_i - \bar{y})}{n}$$

# Propiedades e inconvenientes

La covarianza permite estimar conceptos relativos a la correlación entre las dos variables

I. Su signo indica el sentido de la correlación entre las variables.

- Si  $\text{Cov}_{xy} > 0$ , la correlación es directa.
- Si  $\text{Cov}_{xy} < 0$ , la correlación es inversa.

# Propiedades e inconvenientes

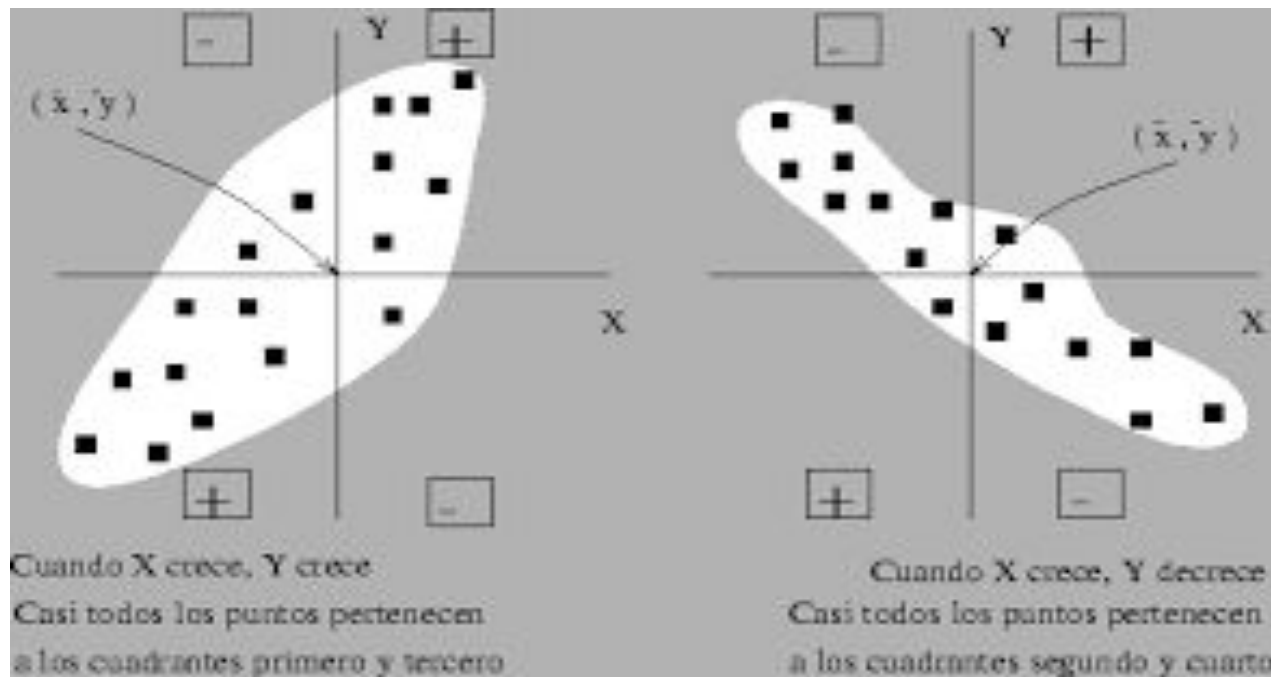
II. Un valor grande de  $\text{Cov}_{xy}$  advierte que la correlación entre las variables puede ser fuerte, pero no lo asegura, no siendo interesante la comparación de dos distribuciones por la covarianza.

Sólo da el sentido de la correlación:  
directa si es positiva e inversa si es negativo.

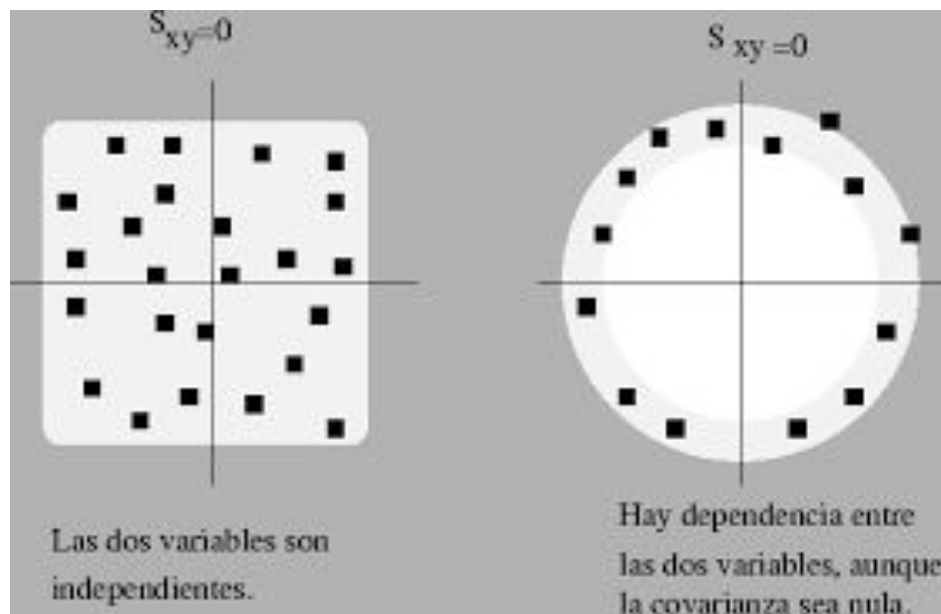




# Gráficamente



# Mas gráficos



## Correlación lineal simple – coeficientes de asociación

Los coeficientes de asociación son valores numéricos que permiten cuantificar el grado de ajuste y de relación lineal entre dos variables.



# Correlación lineal simple – coeficientes de asociación

- **PEARSON**

- Es un coeficiente paramétrico por lo que exige que la m.a. Provenga de una distribución normal
- Solo puede calcularse en variables cuantitativas normales

- **SPEARMAN y KENDALL**

- Son coeficientes no paramétricos

# Pearson

Sea  $(X, Y)$  una variable aleatoria bidimensional. Definimos el coeficiente de correlación lineal entre  $X$  e  $Y$  como

$$r = r_{xy} = \frac{\text{Cov}(x, y)}{S_x \times S_y}$$

Este coeficiente nos da una idea del grado de asociación entre las variables aleatorias  $X$  e  $Y$

# Interpretación

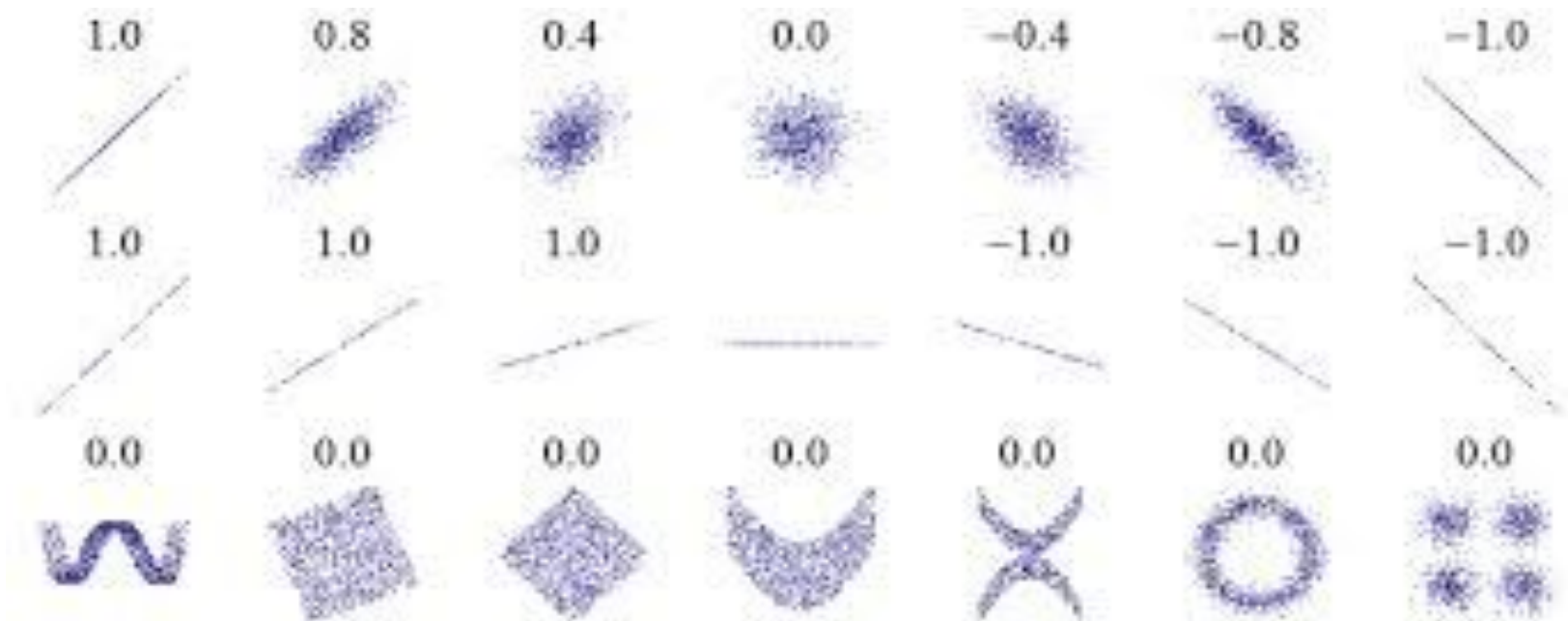
El valor del índice de correlación varía en el intervalo  $[-1,1]$ , indicando el signo el sentido de la relación:

- Si  $r = 1$ , existe una correlación positiva perfecta. El índice indica una dependencia total entre las dos variables denominada *relación directa*: cuando una de ellas aumenta, la otra también lo hace en proporción constante.
- Si  $0 < r < 1$ , existe una correlación positiva.

# Interpretación

- Si  $r = 0$ , no existe relación lineal. Pero esto no necesariamente implica que las variables son independientes: pueden existir todavía relaciones no lineales entre las dos variables.
- Si  $-1 < r < 0$ , existe una correlación negativa.
- Si  $r = -1$ , existe una correlación negativa perfecta. El índice indica una dependencia total entre las dos variables llamada *relación inversa*: cuando una de ellas aumenta, la otra disminuye en proporción constante.

# Gráficamente





# Coeficiente de correlación lineal. Spearman

Sea  $(X, Y)$  una variable aleatoria bidimensional. Pero ahora no podemos afirmar que la distribución de la v.a. es normal o tenemos poco datos. Además nuestra variable es del tipo ordinal (una escala de valores). En estos casos podemos utilizar el coeficiente de Spearman.

Este coeficiente analíticamente tiene un cálculo tedioso

# Coeficiente de correlación lineal. Tau de Kendall

Tau de Kendall o como también es llamado, Coeficiente de Correlación por Rangos de Kendall, es una medida de asociación no paramétrica utilizada para estudiar variables cualitativas ordinales o de razón. Estas variables son distribuidas en categorías con varios niveles que cumplen un orden, por ejemplo, muy bajo, bajo, medio, alto y muy alto.



# Coeficiente de correlación lineal. Tau de Kendall

- Sólo se puede aplicar a partir de tablas cuadradas.
- Las variables utilizadas deben ser de nivel ordinal, intervalo o razón Su resultado debe encontrarse en el rango de -1 a 1.
- Tiene sentido su aplicación, si las variables objeto de estudio no poseen una distribución poblacional conjunta normal

# Coeficiente de correlación lineal. Cálculos

Pandas hace que encontrar correlaciones sea extremadamente fácil. Podemos usar el método `corr` para calcular la correlación por pares de columnas usando los métodos de Pearson, Kendall o Spearman. En este punto, debe tener una hipótesis bien formulada o al menos una idea de lo que le gustaría probar o refutar con la ayuda de correlaciones y saber en que situaciones que coeficientes

# Coeficiente de correlación lineal. Spearman

```
import pandas as pd
```

```
df = pd.read_csv("midataset.csv")
```

```
df.corr(method='pearson')
```

```
df.corr(method='spearman')
```

```
df.corr(method='kendall')
```

# Coeficiente de correlación lineal. Spearman

```
import pandas as pd
```

```
df = pd.read_csv("midataset.csv")
```

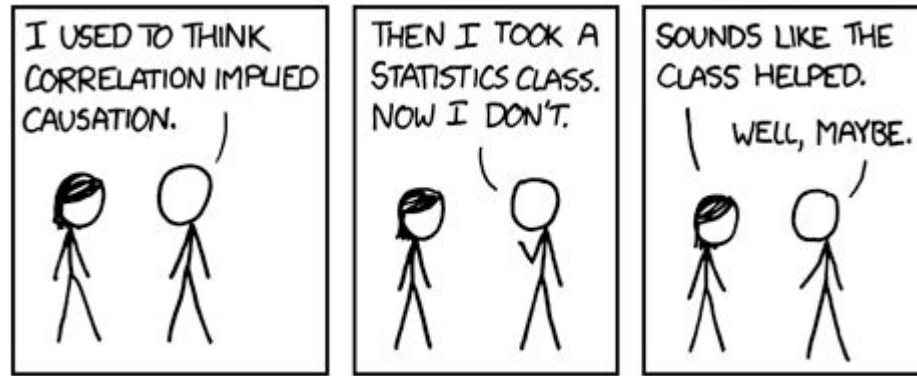
```
df.corr(method='pearson')
```

```
df.corr(method='spearman')
```

```
df.corr(method='kendall')
```

# No todo es magia

Recuerden que la correlación no implica causalidad. Por ejemplo, si las ventas de helados están correlacionadas positivamente con los ataques de los tiburones a los nadadores, eso no significa que el consumo de helados de alguna manera hace que los tiburones ataquen. Otra variable, como el clima cálido, puede provocar un aumento tanto en las ventas de helados como en las visitas a las playas.





¿Dudas?

