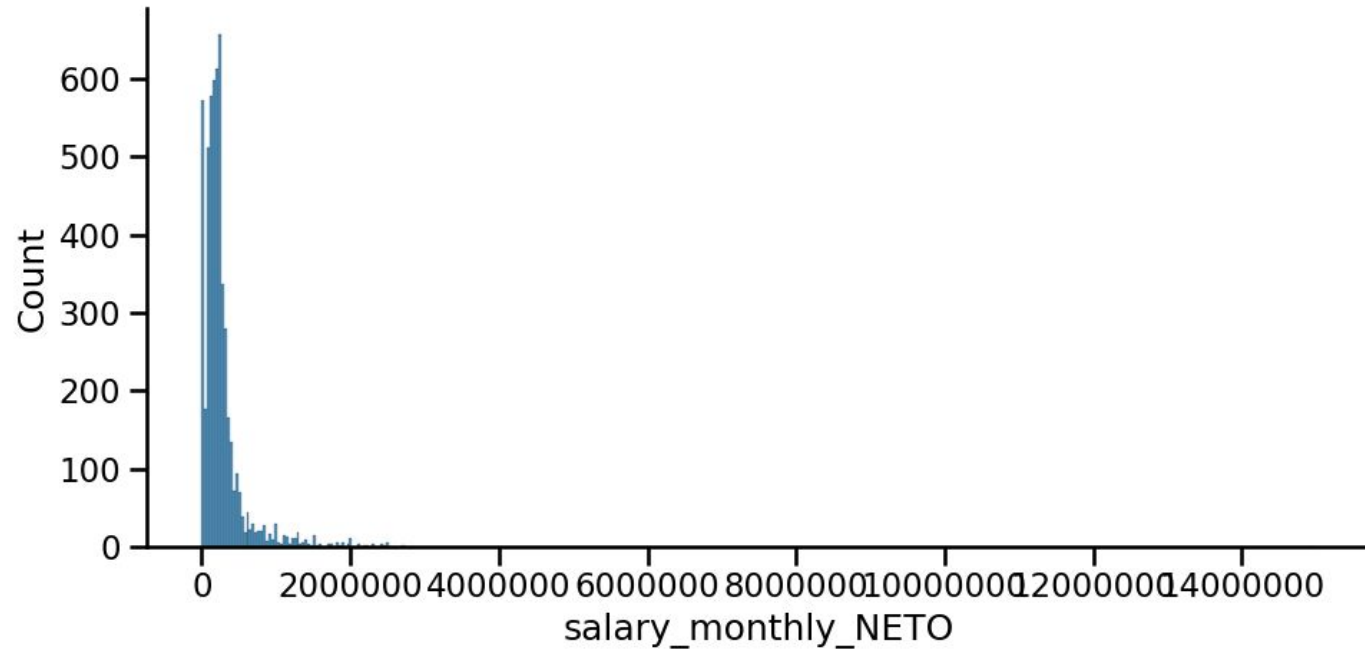


Análisis y Visualización de Datos



Diplomatura CDAAyA 2024



¿Más herramientas para describir el “comportamiento” de los datos de la columna del sueldo?

Estadística Descriptiva

Estadística Descriptiva

Representa la información de una manera distinta para facilitar su interpretación, pero no permite realizar predicciones o inferencias

Análisis de frecuencias

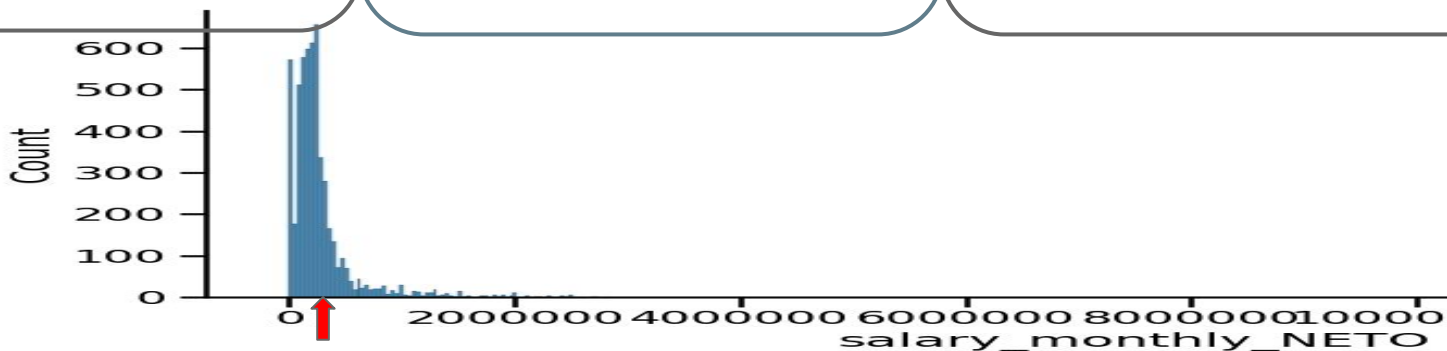
¿Cuánto ocurre cada uno de los valores (o en intervalos) de un conj de datos?

Medidas de tendencia central

¿Cuál es el valor más representativo del conj. de datos?

Medidas de dispersión

¿Qué tan alejados están los datos de la tendencia central?



Medidas de tendencia central

una muestra de
datos numéricos

$$x_1, x_2, \dots, x_N$$

(de una columna
en un DataFrame)

##	datos
## 1	40
## 2	60
## 3	50
## 4	45
## 5	65
## 6	70
## 7	95
## 8	90
## 9	45
## 10	60
## 11	43
## 12	56
## 13	65
## 14	80
## 15	45
## 16	70
## 17	45
## 18	75
## 19	45
## 20	54
## 21	35
## 22	46
## 23	47
## 24	50
## 25	50
## 26	60
## 27	50
## 28	50

La **media muestral** (aritmética) o promedio se
calcula como:

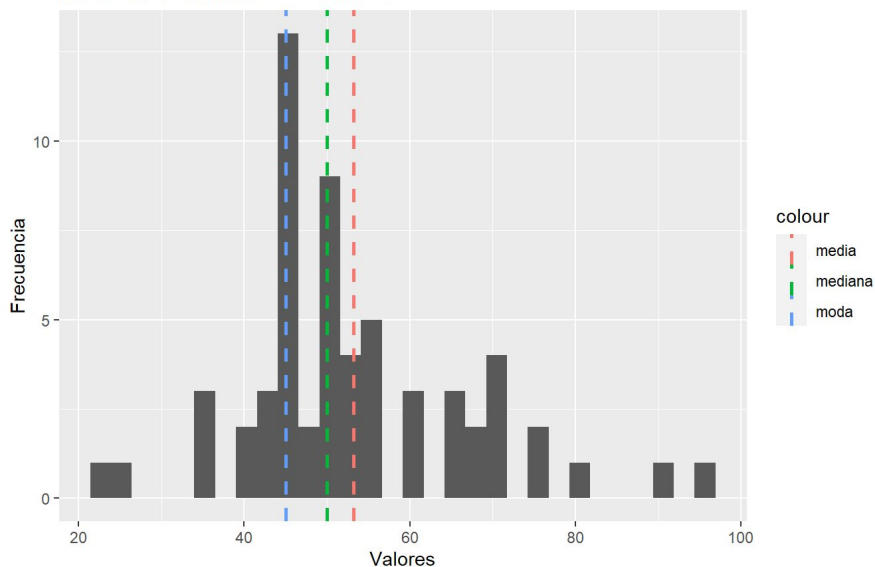
$$\bar{x} = \frac{1}{N} \sum_i^N x_i$$

Medidas de tendencia central

x_1, x_2, \dots, x_N numéricos

Histograma de los datos

Media= 53.17 Mediana = 50 Moda= 45



La **mediana** se calcula como:

1. Ordenar las realizaciones de menor a mayor
2. Si N es impar, la mediana es el valor central:

$$\text{mediana} = x_{(N+1)/2}$$

1. Si N es par, la mediana es el promedio de los dos valores centrales:

$$\text{mediana} = (x_{N/2} + x_{N/2+1})/2$$

La **moda** es el dato con mayor frecuencia, para

numéricas **intervalo modal** depende del histograma

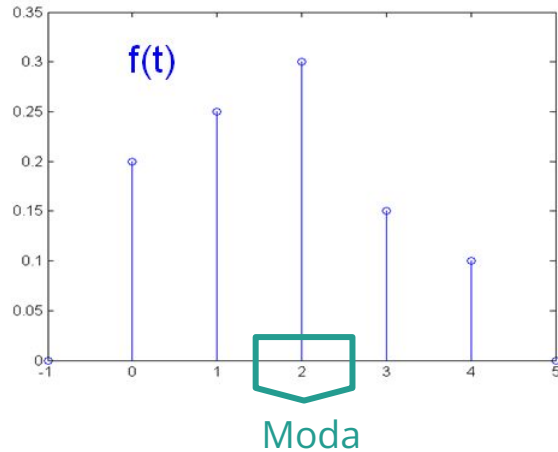
Medidas de tendencia central

si ahora

x_1, x_2, \dots, x_N

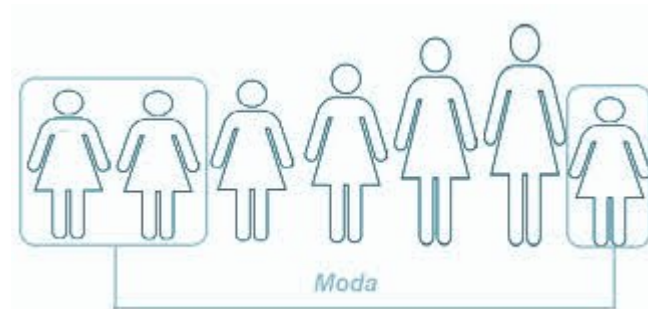
categoricos

p/ $x_i = X_i(\omega) \in \{c_1, \dots, c_k\}$



La **moda** es el dato con mayor frecuencia, el que más se repite

Sólo hay más de una moda cuando el conteo de dos valores es igual.

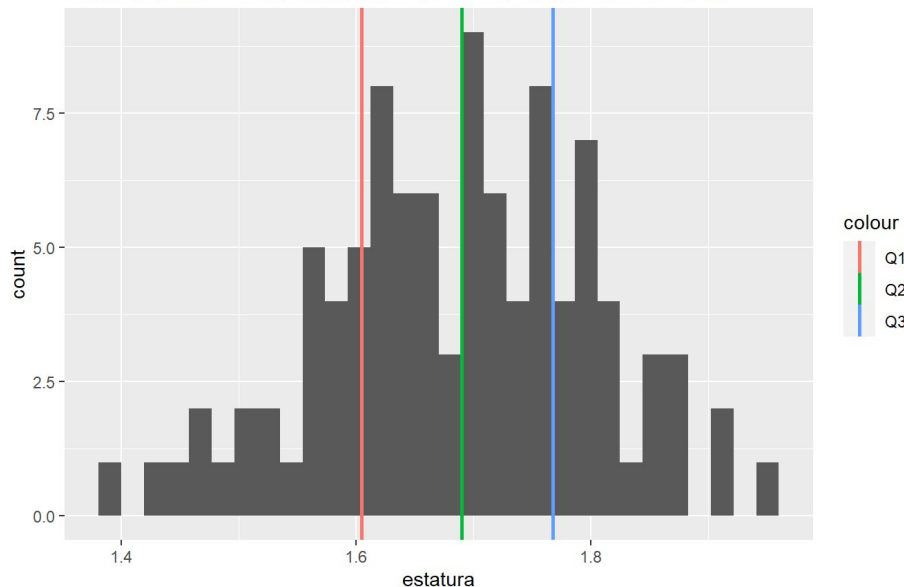


Medidas de posición

x_1, x_2, \dots, x_N numéricos

Histograma de los estatura

Cuartil 1 al 25% = 1.61 , Cuartil 2 al 50% = 1.69 , Cuartil 3 al 75% = 1.77



El **percentil-k** es el valor x_i tal que el $k\%$ de los valores de la muestra son menores a x_i .

No hay una única fórmula para calcular los percentiles, pero en general:

1. Ordenar las realizaciones tal que $x_i \leq x_{i+1}$
2. Seleccionar el elemento de la serie en la posición: menor entero mayor o igual a $k \cdot N / 100$.

percentil 25 es el primer **cuartil** Q1

percentil 50 el segundo **cuartil** Q2 (mediana)

percentil 75 el tercer **cuartil** Q3

Medidas de dispersión

x_1, x_2, \dots, x_N numéricos

La **varianza muestral** mide la variación de los datos a través de la distancia cuadrada a la media muestral.

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

La **desviación estándar** es la raíz cuadrada de la varianza. Está en la misma unidad que los datos.

El **coeficiente de variación** es la desviación estándar dividida la media muestral. Es comparable entre distintas v.a.

Medidas de dispersión

Datos numéricos

X_1, X_2, \dots, X_N

⋮

El rango y el rango intercuartílico miden en qué intervalo se encuentran un cierto porcentaje de los datos.

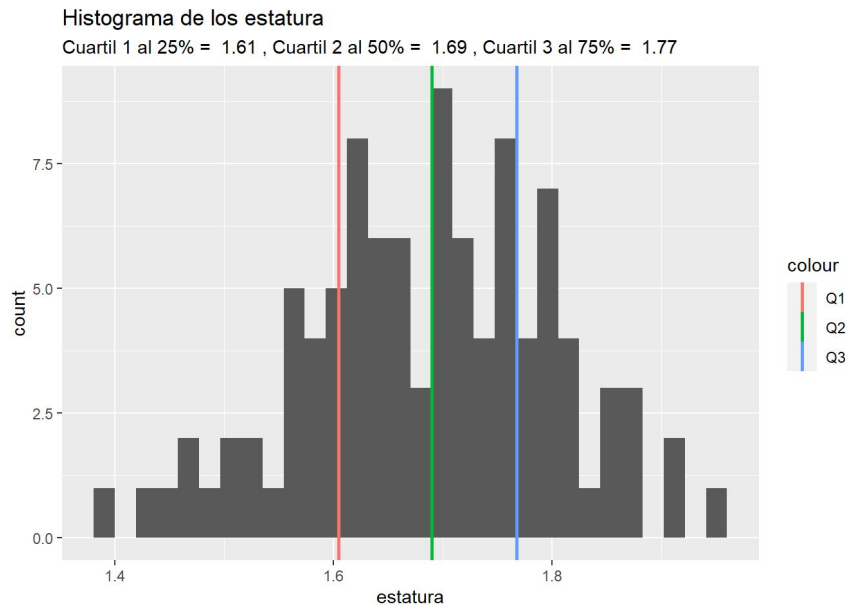
Rango:

percentil-100 - percentil-0

Rango intercuartílico:

percentil-75 - percentil-25

$Q3 - Q1$



Usos de los percentiles y rangos

- En el caso de la mediana (percentil-50), medir la tendencia central
- Contextualizar el valor de un dato con respecto a otros

Una persona de sexo femenino de 6 años mide 95cm

Está en el 10% de personas con menor estatura del mismo grupo.

[[Curvas](#)], se visualiza, para cada edad, los percentiles de la distribución condicional al grupo.

- Identificación y eliminación de valores extremos

Demo con Notebook

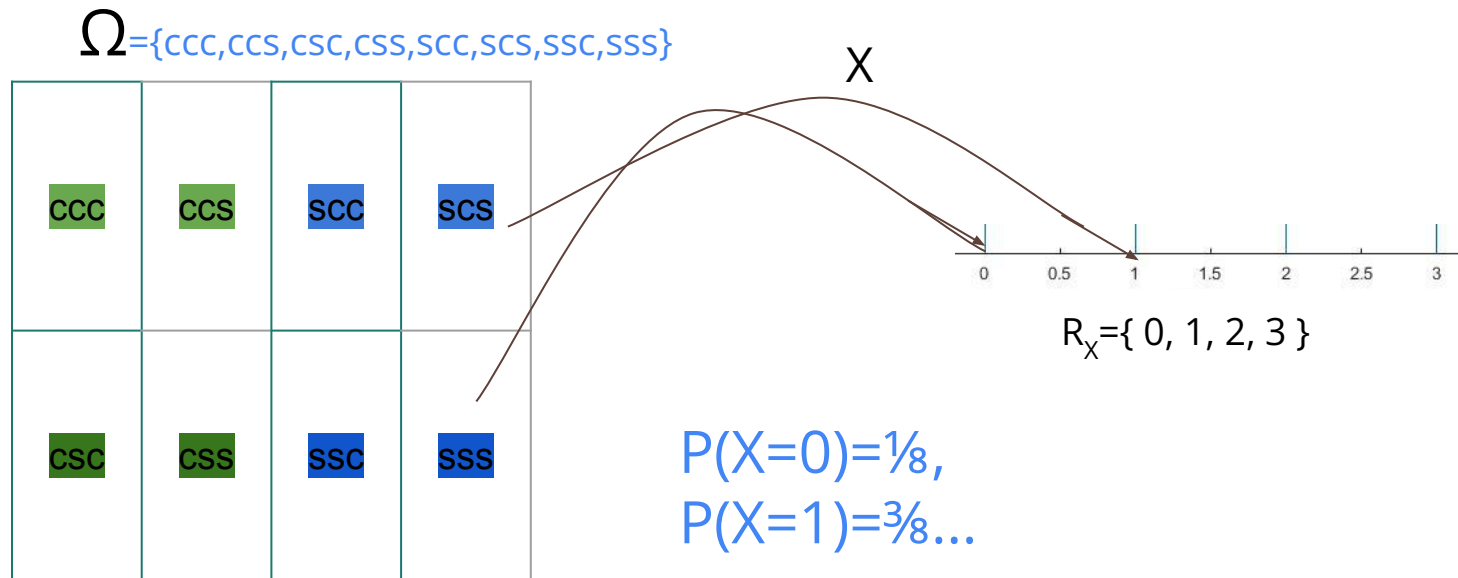
02 Datos y Modelos.ipynb

**Teoría, datos, experimentos,
simulación...**

**¿Que es todo esto y cómo se
combinan?**

Variable Aleatoria (discreta numérica)

X = cantidad de caras en 3 tiradas de moneda.



Variable Aleatoria (repetición del experimento)

X= cantidad de caras en 3 tiradas de moneda.

```
C, S = 'c', 's'
SAMPLE_SPACE = ['-'.join(x) for x in
                 itertools.product([C, S], repeat=3)]
SAMPLE_SPACE

['c-c-c', 'c-c-s', 'c-s-c', 'c-s-s', 's-c-c', 's-c-s', 's-s-c', 's-s-s']
```

```
sampled_values = [
    x.count(C) for x in numpy.random.choice(SAMPLE_SPACE, 1000)]
```

Proporción de resultados tal que $X=k$:

```
result = numpy.unique(sampled_values, return_counts=True)
[(label, count/1000.0) for label, count in zip(*result)]

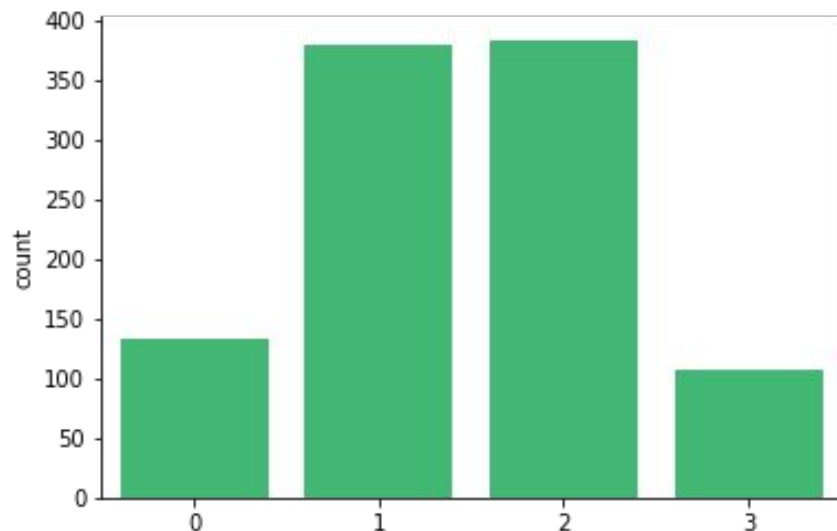
[(0, 0.132), (1, 0.379), (2, 0.383), (3, 0.106)]
```

Variable Aleatoria (repetición del experimento)

```
result = numpy.unique(sampled_values, return_counts=True)  
[(label, count/1000.0) for label, count in zip(*result)]
```

```
[(0, 0.132), (1, 0.379), (2, 0.383), (3, 0.106)]
```

la Proporción de la muestra tal
que $X=k$, estima la probabilidad
 $P(X=k)$, $p/ k=0,1,2,3$



Variable Aleatoria (modelo matemático)

X = cantidad de caras en 3 tiradas de moneda. $p(k) = P(X=k)$?

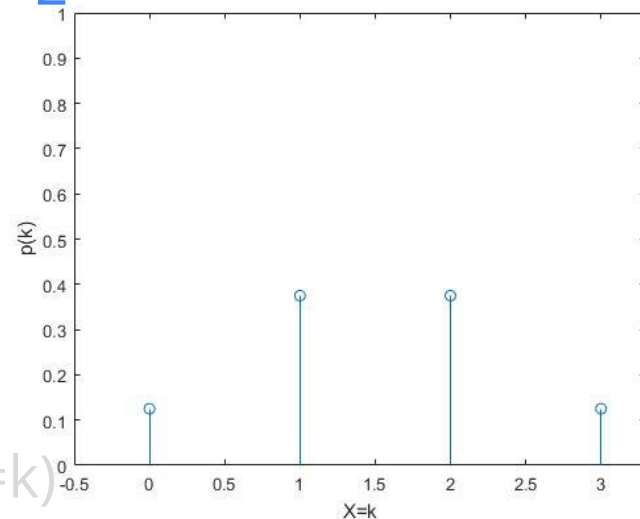
$$\Omega = \{ccc, ccs, csc, css, scc, scs, ssc, sss\}, \quad \#\Omega = 8 = 2^3$$

$$p(0) = P(X=0) = 1/8$$

$$p(1) = P(X=1) = 3/8$$

$$p(2) = P(X=2) = 3/8$$

$$p(3) = P(X=3) = 1/8$$



Notar que la suma da 1, $\sum_k p(k) = 1 = \sum_k P(X=k)$

Variable binomial

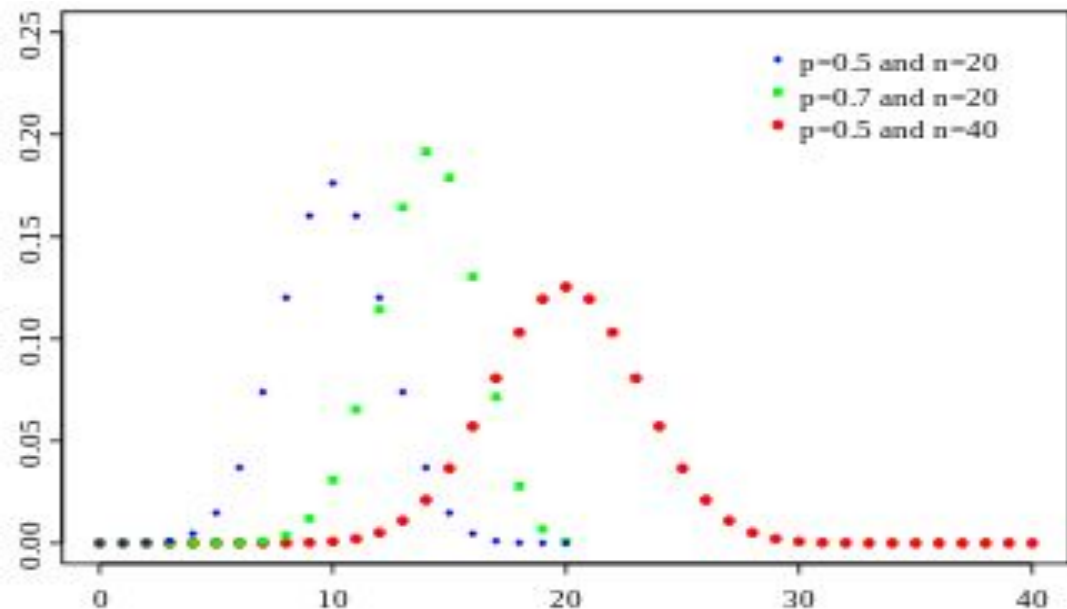
Sea X la v. a. discreta modela: cantidad de “éxitos” en una n -upla

$$P(X=k) = \frac{n!}{(n-k)! k!} p^k (1-p)^{(n-k)}$$

$k=0,1,\dots,n$

p =probabilidad de “éxito”.

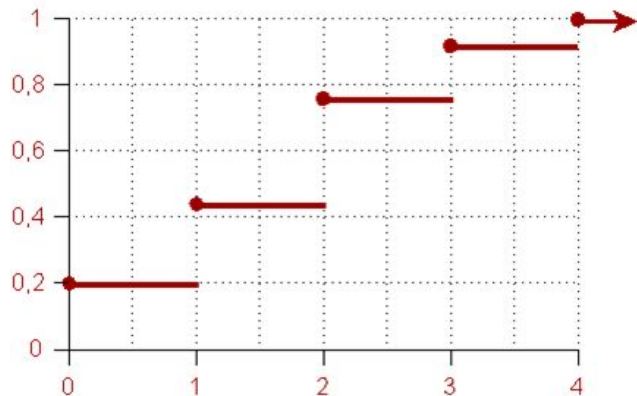
$X \sim B(n,p)$, ejemplos?



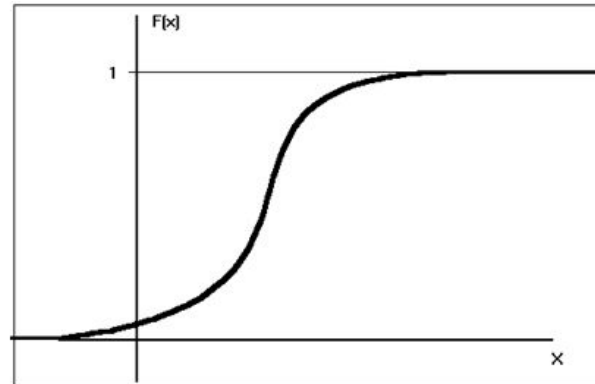
Función de Distribución Acumulada

La Función de Distribución Acumulada de la v.a. X , es la función $F: \mathbb{R} \rightarrow [0,1]$ definida por

$$F(t) = P(X \leq t) = P(\{\omega / X(\omega) \leq t\})$$



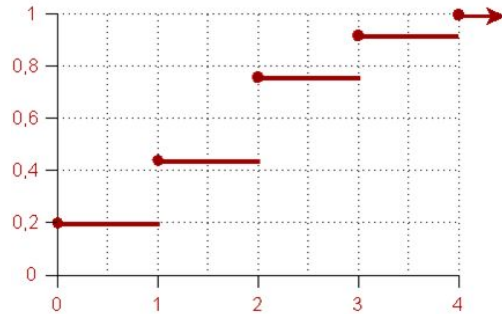
X discreta



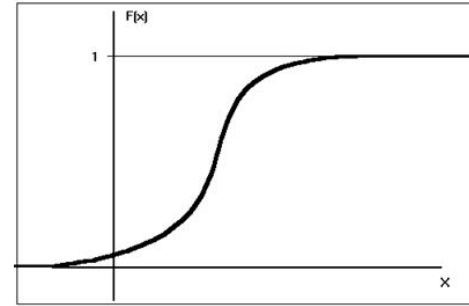
X continua

Función de densidad

FDA: $F(t) = P(X \leq t)$



X discreta

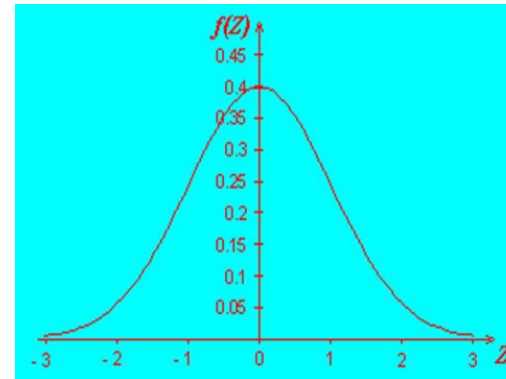
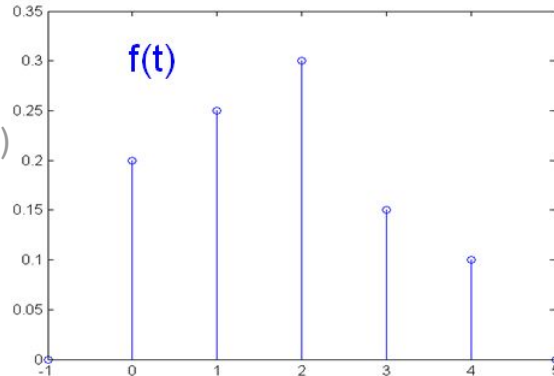


FDA: $F(t) = P(X \leq t)$

X continua

densidad discreta
(probabilidad puntual)

$f(t) = P(X=t)$



densidad:

$f(t) = F'(t)$

$F(t) = \int^t f(x) dx$

Propiedades de función de densidad

1) $f(t) \geq 0$ para todo t

2) $\int f(t) dt = 1$ para variables continuas y (entre $-\infty$ y $+\infty$)

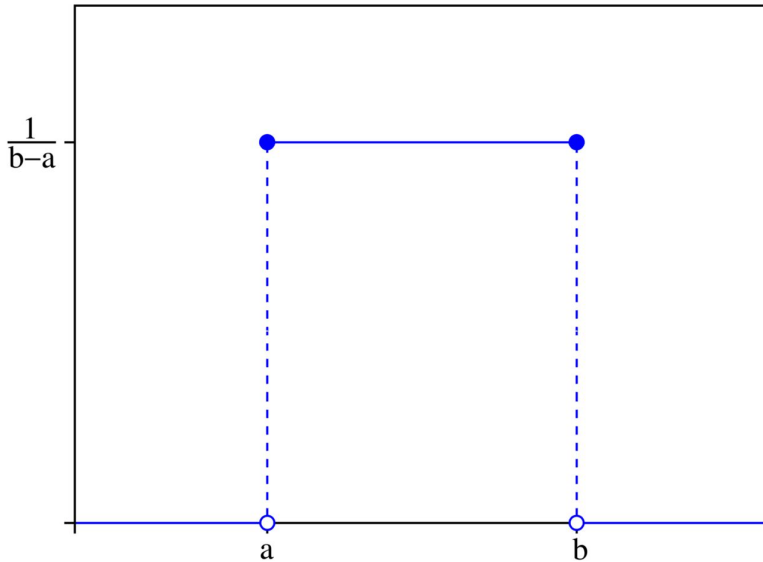
2) $\sum f(t) = 1$ para variables discretas (para todos los valores)

cualquier función que cumple con 1 y 2 es una función de densidad de alguna v. a.

Distribución Uniforme

X v.a. tiene **distribución uniforme** si su función **densidad** es

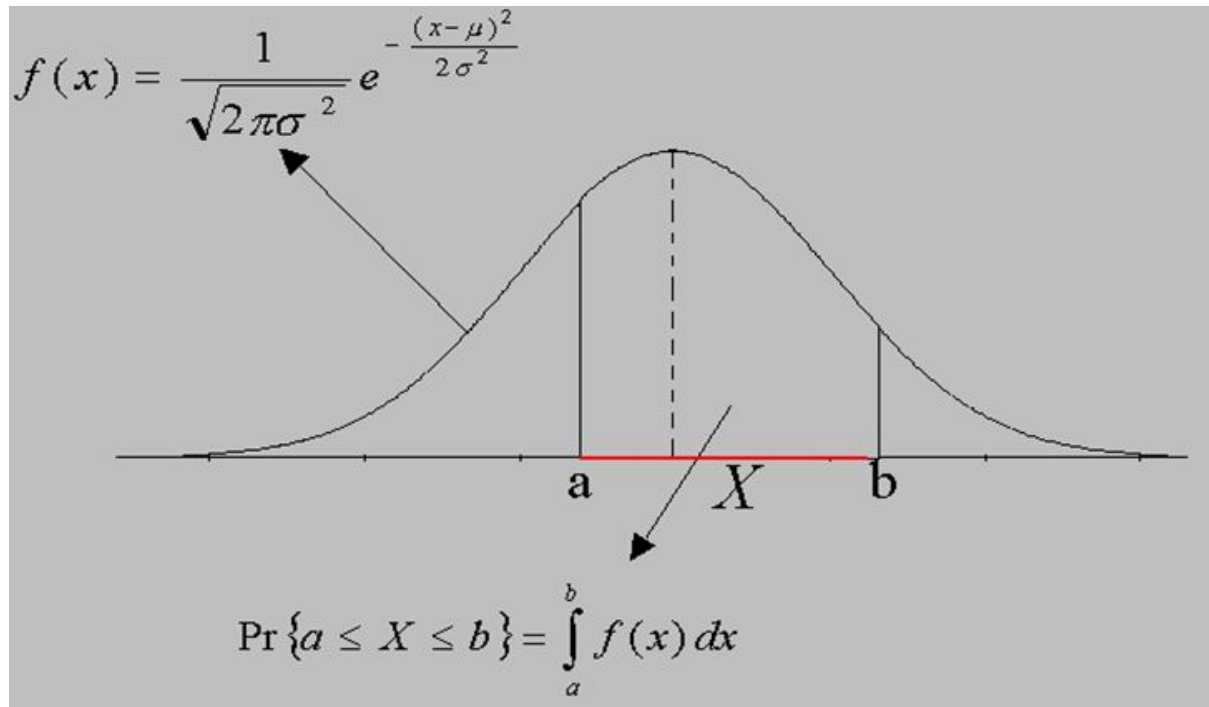
$$f(t) = 1/(b-a) \text{ si } a \leq t \leq b, 0 \text{ c.c.}$$



Notación $X \sim U(a,b)$, $a < b$ parámetros

Distribución Normal o Gaussiana

X v.a. continua tiene distribución normal (Gaussiana) si su función de densidad es la siguiente:



Con $\mu \in \mathbb{R}$ y $\sigma^2 \in (0, \infty)$

parámetros

Notación $X \sim N(\mu, \sigma^2)$

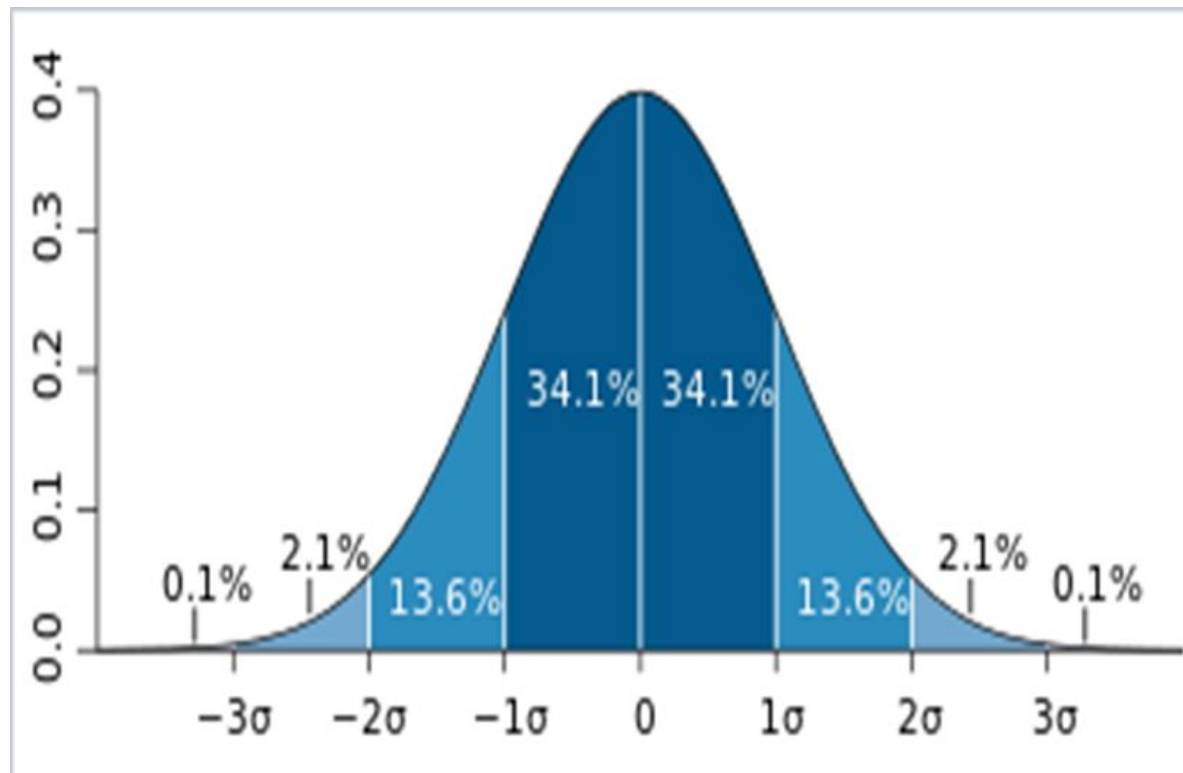
Distribución Normal o Gaussiana

$$X \sim N(0, \sigma^2)$$

si además $\sigma^2=1$

$X \sim N(0,1)$, se dice

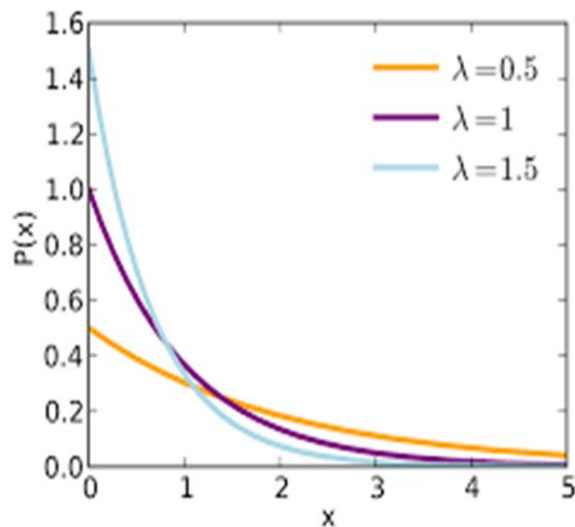
Normal Estándar



Distribución Exponencial (caso especial de Gamma)

X v.a. tiene distribución exponencial si su densidad es:

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{si } x \geq 0 \\ 0 & \text{si } x < 0 \end{cases}$$



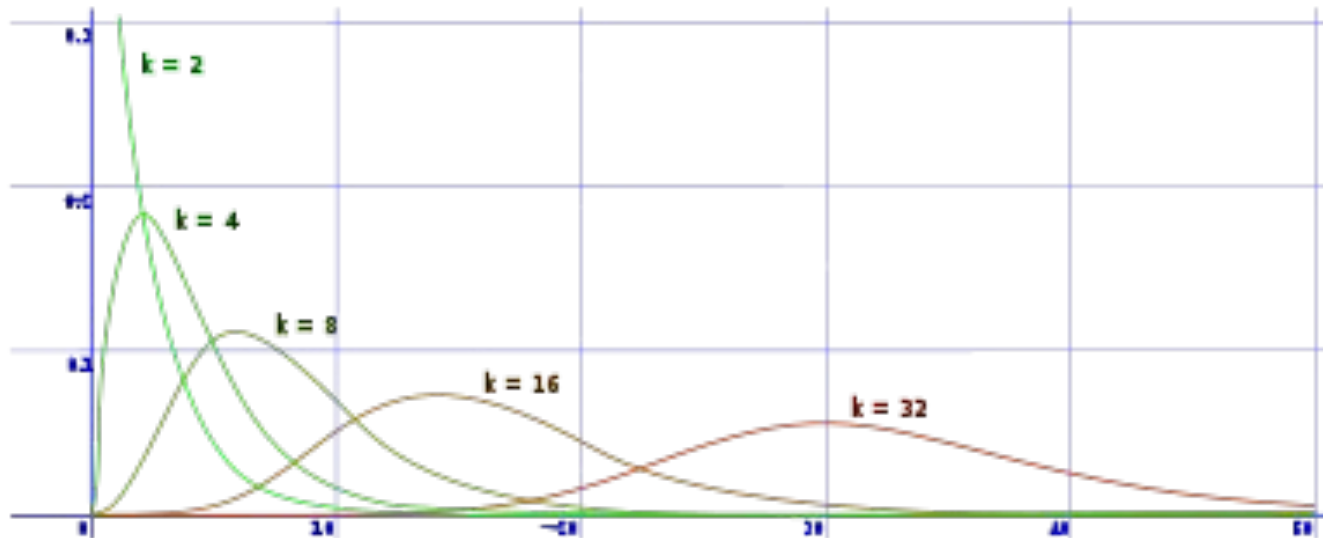
Notación $X \sim \text{Exp}(\lambda)$, $\lambda > 0$ parámetro

suele utilizarse para modelar tiempo de espera

Distribución Chi Cuadrado

Diremos la v.a. X tiene distribución Chi- cuadrado con k grados de libertad.

Notación $X \sim \chi_k^2$ si su función de densidad está dada por:



Medidas estadísticas de una v.a. o de una densidad

X v.a. numérica con densidad f

- **Media o Esperanza** de X (Medida de posición):

$$\mu = E(X) = \int t f(t) dt \quad \text{ó} \quad \mu = E(X) = \sum t f(t), \text{promedio ponderado por la densidad } (\mu \in \mathbb{R})$$

- **Varianza** (Medidas de dispersión):

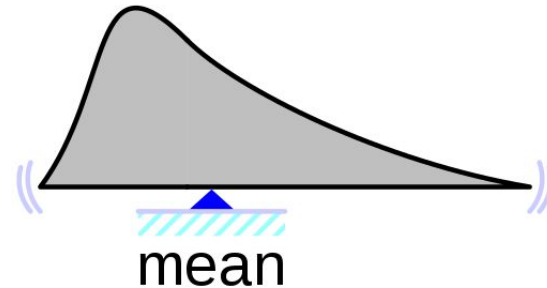
$$\sigma^2 = \text{Var}(X) = E((X-\mu)^2) = \int (t-\mu)^2 f(t) dt \quad \text{ó} \quad \sigma^2 = E((X-\mu)^2) = \sum (t-\mu)^2 f(t) \quad (\sigma^2 \in \mathbb{R}^+)$$

En una va con densidad normal coinciden con los parámetros μ y σ^2 respectivamente

Media

Media Muestral $\sum_{i=1}^n x_i / n$, (promedio) vs

Media o Esperanza de una v.a. X , $\mu = E(X) = \int t f(t) dt$ ó $\mu = E(X) = \sum t f(t)$

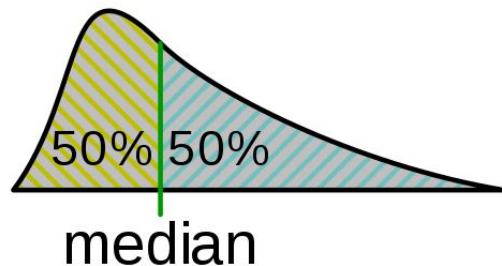


Mediana

Se ordena la muestra de menor a mayor: $x_{(1)}, \dots, x_{(n)}$ y se calcula...

Mediana Muestral vs

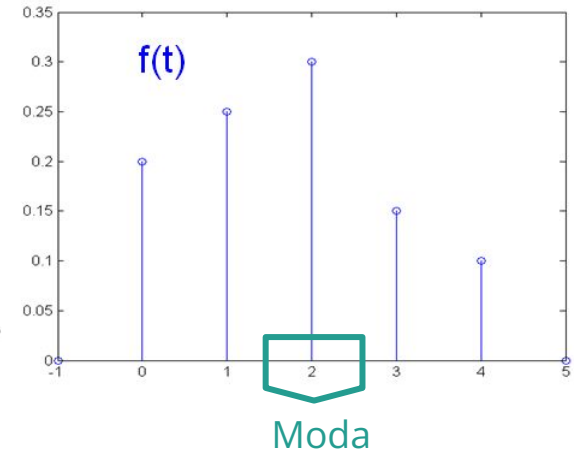
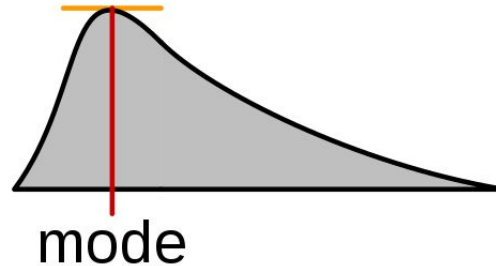
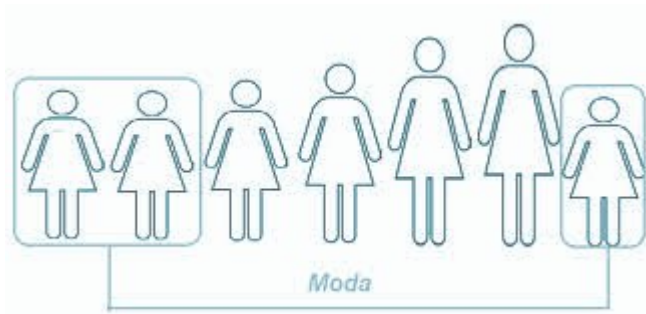
Mediana de una v.a. X , o de su densidad es x_e tal que $P(X \leq x_e) = P(X \geq x_e)$



Moda

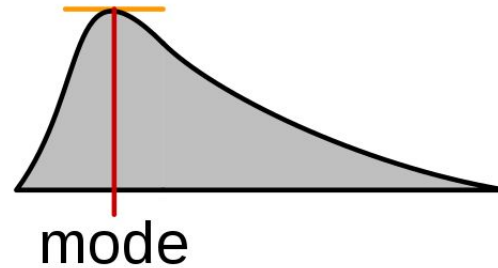
Resultado (o intervalo) con mayor frecuencia en la **muestra**. vs

Valor con **mayor probabilidad** o **densidad** x_0 tal que $f(x_0) \geq f(x)$, p/ todo x

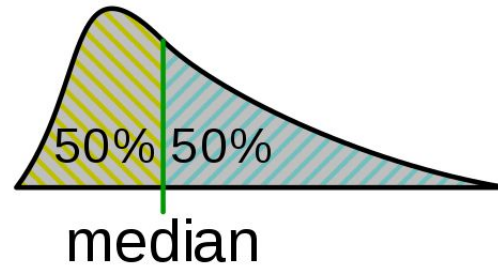


Comparación de Medidas

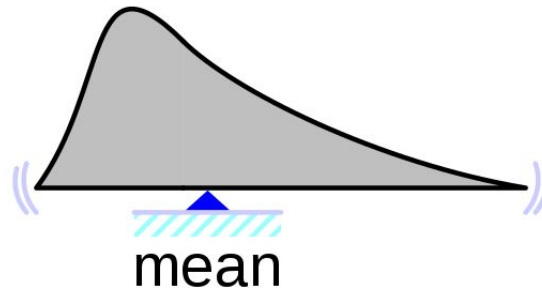
Moda:



Mediana:



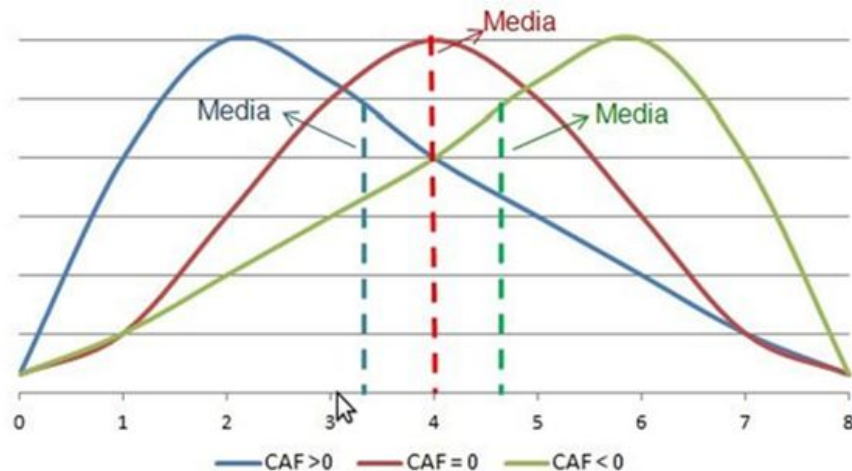
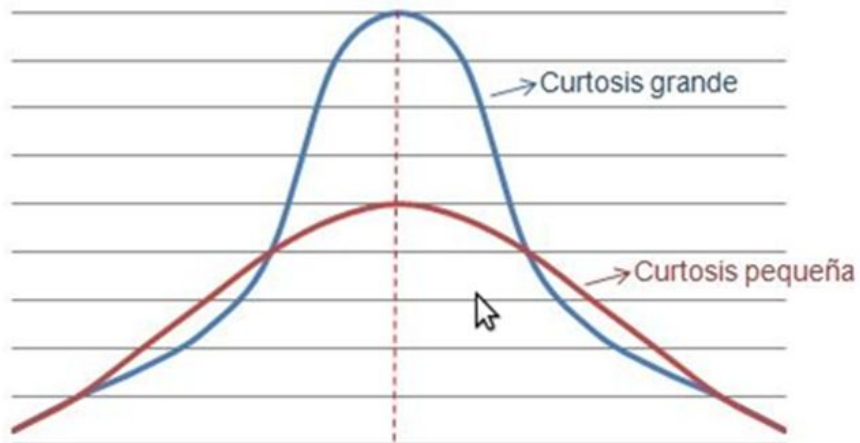
Media:



Otras Medidas, del modelo (de una v.a.)

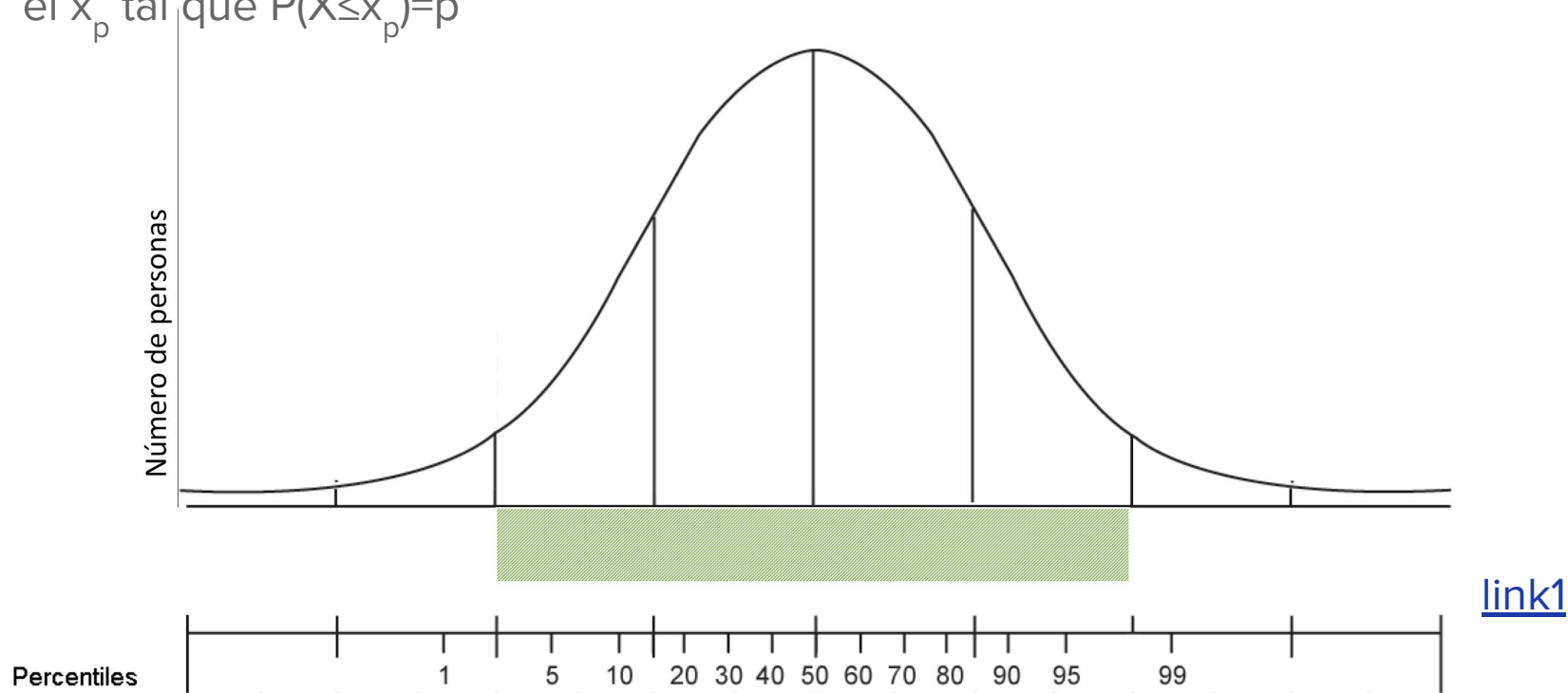
Dada una función de densidad f (de una v.a. X) se define:

Desvío: $\sigma = (\sigma^2)^{1/2} = (\text{Var}(X))^{1/2}$ - **Kurtosis:** $E((X-\mu)^4)/\sigma^4$ - **Sesgo/Asimetría:** $E(X-\mu)^3/\sigma^3$



Percentiles

El percentil es una medida de posición. El p -ésimo percentil o percentil $p \times 100\%$, es el x_p tal que $P(X \leq x_p) = p$



Algunas propiedades de v.a. (modelo) y su distribución

- Si $X \sim N(\mu, \sigma^2)$ y $Z = (X - \mu)/\sigma$, entonces $Z \sim N(0, 1)$
- Si $Z \sim N(0, 1)$, entonces $Z^2 \sim \chi_1^2$ Chi cuadrado con 1 gl

Datos vs modelos

Medidas a partir de datos ↔ Medidas muestrales

Sean los n datos de una muestra: x_1, \dots, x_n (observaciones de la v.a.)

Media muestral (promedio): $x_M = \sum_{i=1}^n x_i / n = \bar{X}$

Varianza muestral : $\sum_{i=1}^n (x_i - x_M)^2 / n$

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

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

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

Tendencia

La tendencia habitual, si los datos están descritos en los términos de \bar{X} y S_X (desvío), es hacer aquellas típicas inferencias que **sólo son ciertas si la distribución de los datos se ajusta bien a la distribución normal**:

- $\bar{X} \pm S_X$ supone el 68.5% aproximadamente de la población,
- $\bar{X} \pm 2S_X$ supone el 95% aproximadamente de la población
- $\bar{X} \pm 3S_X$ supone el 99.5% aproximadamente de la población

Bondad de ajuste

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

- Gráficos QQ (Quantil muestral vs Quantil modelo)

Dentro de los test más usados para normalidad:

- Test de Kolmogorov-Smirnov (Test KS)

(En próxima semana veremos Test de Hipótesis)

En una frase:
¿Cuánto cobran
l@s
programadores
en Argentina?

¿Respuestas?

**¿Qué pregunta
respondimos en
realidad?**

**¿Cuánto cobran l@s
programadores
experimentados en
Argentina?**

¿Afecta el nivel de estudios en el salario de l@s programador@s en Argentina? ¿Cómo?

Ejercicio

Seguir el proceso de análisis
propuesto:

1. Hipótesis
2. Análisis de v.a.
3. Experimento