

# Datos

- Descripción de una sola variable.
  - Datos y distribuciones de frecuencias.
  - Medidas de centralización y dispersión.
  - Medidas no centrales.
  - Visualización datos con medidas no centrales
  - Medidas de asimetría y curtosis.
  - Datos atípicos y diagramas caja (Boxplots).
- Descripción conjunta de varias variables.
  - Distribuciones de frecuencia multivariantes
  - Medidas de dependencia lineal
  - Recta de Regresión
  - Matriz de varianzas

# Datos: Descripción de una sola variable

## Descripción de una sola variable:

- Datos y distribuciones de frecuencias
- Representaciones gráficas
- Medidas de centralización y dispersión
- Medidas no centrales
- Visualización datos con medidas no centrales
- Medidas de asimetría y curtosis
- Datos atípicos y diagramas caja (boxplots)

# Datos y distribuciones de frecuencias

- Recordemos que la estadística descriptiva estudia los procedimientos para sintetizar información para un conjunto de datos de una variable  $x$ :
  - Variables cualitativas, categóricas o atributos.
  - Variables cuantitativas discretas.
  - Variables cuantitativas continuas.
- **Frecuencia absoluta** de un suceso  $x_i$ :
  - Es el número de veces que se observa  $x_i$
- **Frecuencia relativa** de un suceso  $x_i$  :
  - $f_r(x_i) = \# \text{ veces que se observa } x_i / \# \text{ total de datos}$
  - Cumple que  $\sum f_r(x_i) = 1$

Parte de la foto sacada de

[https://es.wikipedia.org/wiki/Frecuencia\\_estad%C3%ADstica](https://es.wikipedia.org/wiki/Frecuencia_estad%C3%ADstica)

A	B		
$x_i$	$n_i$	$f_i$	$p_i$
1	16	8/25	32%
2	20	2/5	40%
3	9	9/50	18%
4	5	1/10	10%
		$N=50$	$\sum f_i=1$ $\sum p_i=100\%$

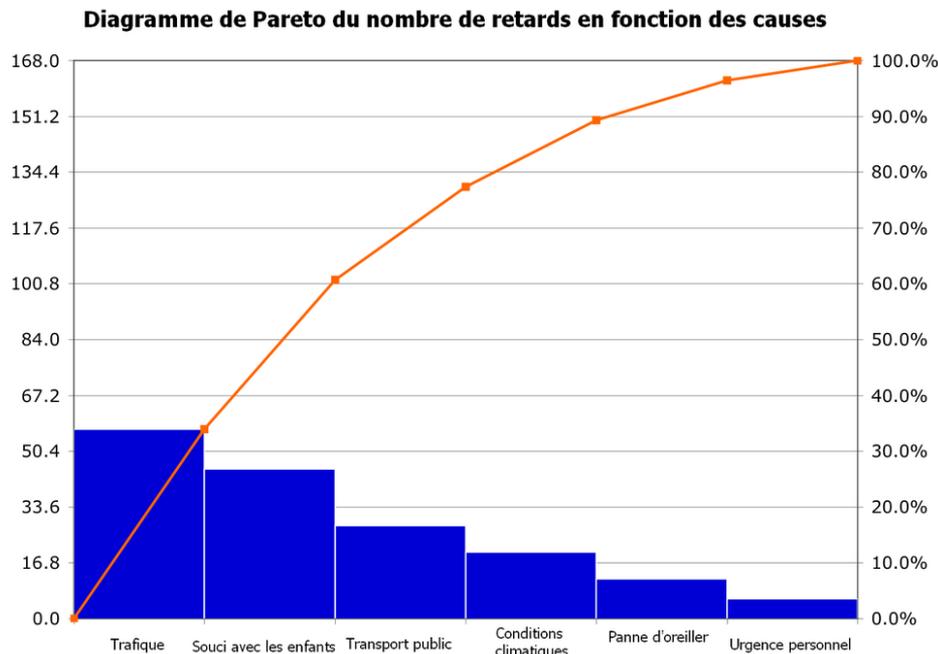
# Representaciones gráficas

- Hay muchas:
  - Diagrama de Pareto.
  - Diagrama de barras.
  - Histogramas.
  - Gráficos temporales.
  - Otras representaciones.

# Representaciones gráficas: Diagrama Pareto

- Para datos cualitativos y representa el principio de Pareto: pocos vitales, muchos triviales.

Foto sacada de [https://es.wikipedia.org/wiki/Diagrama\\_de\\_Pareto#/media/File:Pareto.png](https://es.wikipedia.org/wiki/Diagrama_de_Pareto#/media/File:Pareto.png)

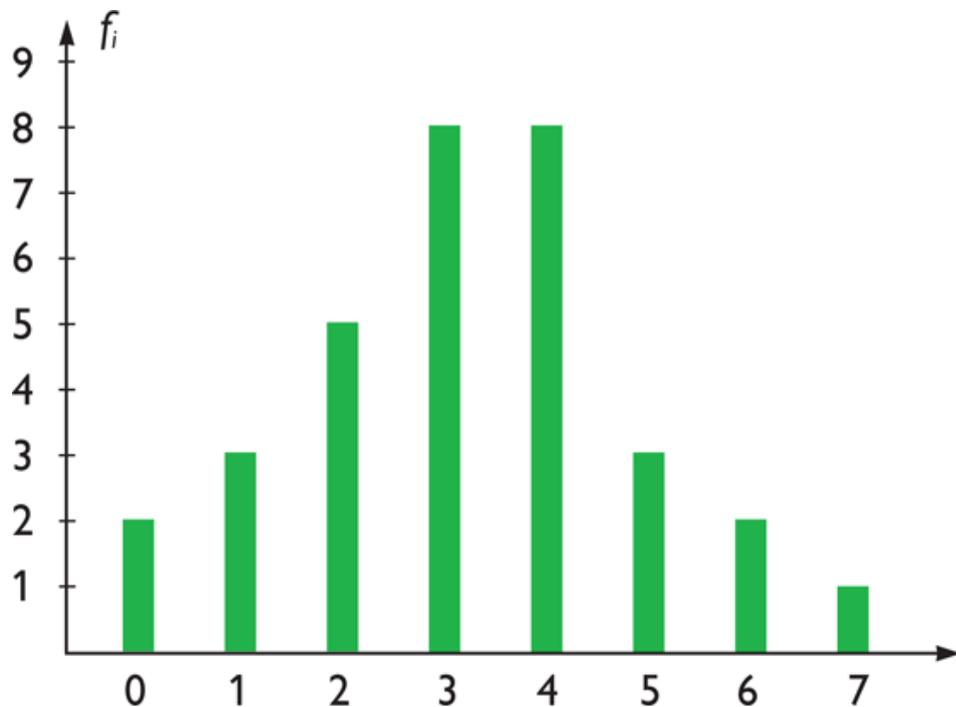


- Organiza los datos de manera que el orden sea descendiente.
- De esta forma se puede asignar orden a las prioridades en el que se toman las decisiones.
- Muestra el principio de Pareto: muchos problemas sin importancia frente a unos pocos muy importantes.

# Representaciones gráficas: Diagrama de barras

- Generalmente para variables discretas.

Foto sacada de  
[http://marcalboran.org/wiki/index.php/Graficos\\_estadisticos](http://marcalboran.org/wiki/index.php/Graficos_estadisticos)



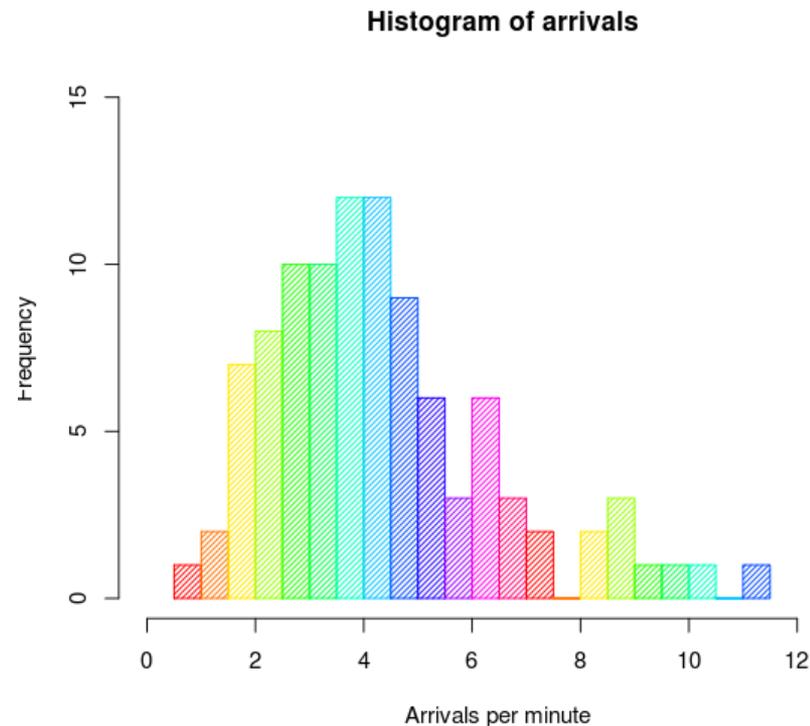
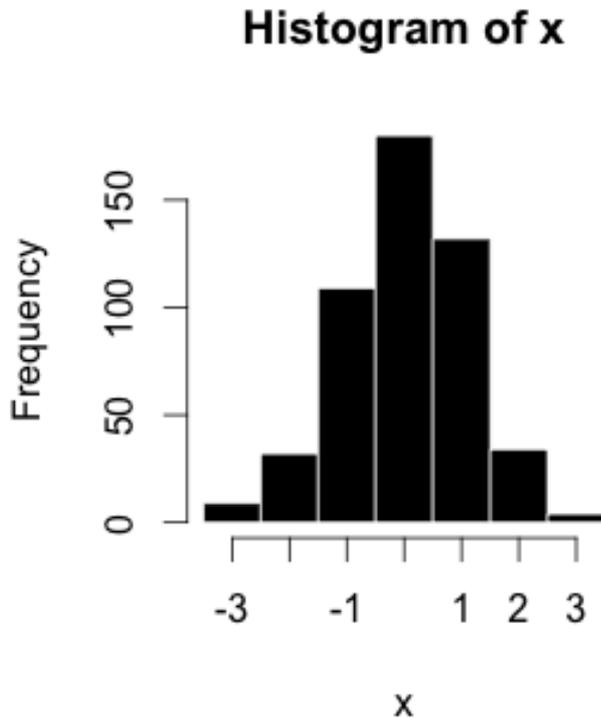
# Representaciones gráficas: Histograma

- Conjunto de rectángulos cada uno de los cuales representa un intervalo de agrupación o clase.
- Sus bases son iguales a la amplitud del intervalo, y las alturas se determinan que su área sea proporcional a la frecuencia de cada clase.
- Vale para valores discretos o continuos.
- $k = \text{ceil}[(\max(x) - \min(x))/h]$ ,
  - donde  $k$  es el número de rectángulos,
  - y  $h$  es la anchura del rectángulo,
- ¿Que  $k$  se selecciona? hay muchas reglas en función del número de muestras,  $n$ :
  - $k = \text{Sqr}(n)$
  - Fórmula de Sturges:  $k = \text{ceil}(\log_2 n + 1)$
  - Regla de Rice:  $k = \text{ceil}(2n^{1/3})$
  - Etc.

# Representaciones gráficas: Histograma

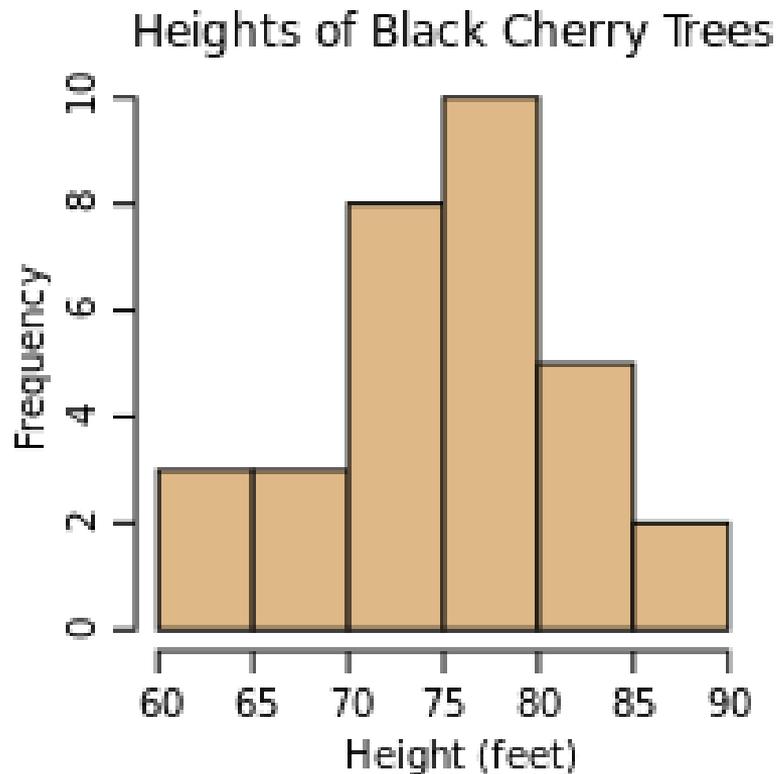
Fotos sacadas de <https://en.wikipedia.org/wiki/Histogram>

Bin	Count
-3.5	9
-2.5	32
-1.5	109
-0.5	180
0.5	132
1.5	34
2.5	4
3.5	9



# Representaciones gráficas: Histograma

- Explicar el significado de un histograma.



Fotos sacadas de

[https://commons.wikimedia.org/wiki/File:Black\\_cherry\\_tree\\_histogram.svg](https://commons.wikimedia.org/wiki/File:Black_cherry_tree_histogram.svg)

	Girth	Height	Volume
1	8.3	70	10.3
2	8.6	65	10.3
3	8.8	63	10.2
4	10.5	72	16.4
5	10.7	81	18.8
6	10.8	83	19.7
7	11.0	66	15.6
8	11.0	75	18.2
9	11.1	80	22.6
10	11.2	75	19.9
11	11.3	79	24.2
12	11.4	76	21.0
13	11.4	76	21.4
14	11.7	69	21.3
15	12.0	75	19.1
16	12.9	74	22.2
17	12.9	85	33.8
18	13.3	86	27.4
19	13.7	71	25.7
20	13.8	64	24.9
21	14.0	78	34.5
22	14.2	80	31.7
23	14.5	74	36.3
24	16.0	72	38.3
25	16.3	77	42.6
26	17.3	81	55.4
27	17.5	82	55.7
28	17.9	80	58.3
29	18.0	80	51.5
30	18.0	80	51.0
31	20.6	87	77.0

# Medidas de centralización y dispersión

## ➤ Medidas de centralización:

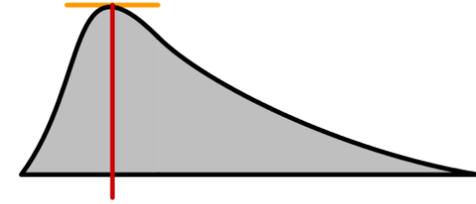
- **Media:**  $\langle x \rangle = (x_1 + x_2 + \dots + x_n) / n = \sum x_i / n$ , o en frecuencias  $\langle x \rangle = \sum x_i \text{fr}(x_i)$ . **Es muy sensible a observaciones atípicas.**
- **Mediana:** valor tal que ordenados en magnitud los datos, el 50% es menor que ella y el 50% mayor. Es el valor central si el número de datos es impar, o la media de los dos centrales si es par. Se suele denominar por **Med.** **No es tan sensible a observaciones atípicas.** Si la media y mediana difieren mucho es posible que exista heterogeneidad de poblaciones.
- **Moda:** valor que se repite más.

## ➤ Medidas de dispersión:

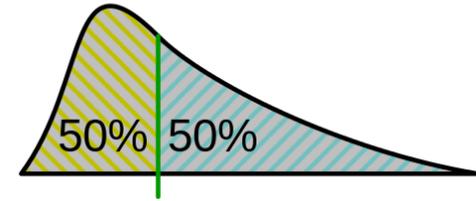
- **Desviación típica:**  $s = [\sum (x_i - \langle x \rangle)^2 / n]^{1/2}$ , o  $s = [\sum (x_i^2 / n) - \langle x \rangle^2]^{1/2}$   
o en frecuencias  $s = [\sum (x_i - \langle x \rangle)^2 f_r(x_i)]^{1/2}$
- **Coefficiente de variación:**  $CV = s / |\langle x \rangle| > 0$  (mide la variabilidad, el inverso es el coeficiente de señal-ruido).
  - Si CV son grandes ( $> 1.5$  normalmente), indican posible error en los datos, o datos heterogéneos.

# Medidas de centralización y dispersión

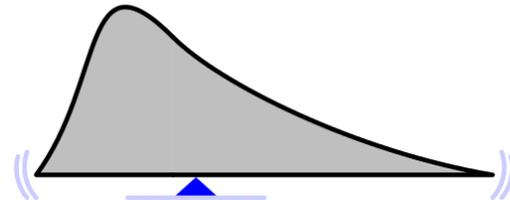
- Medidas de dispersión:
  - Otras medidas de dispersión:
    - Asociadas a mediana:
      - **MEDA** = mediana  $|x_i - \text{Med}|$ ,
      - Como la mediana, MEDA no es tan alterada por datos extremos.
      - Este tipo de medidas que so ven aceptadas por los datos extremos, se denotan como medidas robustas o resistentes.
      - Med y MEDA cumplen que al menos el 50% de los datos esta en el intervalo (Med-MEDA, Med+MEDA).



mode



median



mean

Foto sacada de [https://es.wikipedia.org/wiki/Mediana\\_\(estad%C3%ADstica\)](https://es.wikipedia.org/wiki/Mediana_(estad%C3%ADstica))

# Medidas de centralización y dispersión

- La interpretación que se suele dar a la media y la desviación (información conjunta de media y dispersión):
  - En el intervalo  $(\langle x \rangle - ks, \langle x \rangle + ks)$  existe como mínimo el  $100(1 - (1/k^2))\%$  de las observaciones:
    - Por ejemplo si la media es 500 y la desviación típica 20:
      - Para dos desviaciones típicas tenemos que en el intervalo (460, 540) estarán como mínimo el  $100(1 - (1/2^2))\% = 75\%$  de las observaciones.
      - Para tres desviaciones típicas tenemos que en el intervalo (440, 560) estarán como mínimo el  $100(1 - (1/3^2))\% = 89\%$  de las observaciones.

# Medidas de centralización y dispersión

- Esto es equivalente a que (demostración en el libro).
  - $f_r(|x_i - \langle x \rangle| \geq ks) \geq 1 - (1/k^2)$  que nos permite concluir que en cualquier distribución de datos se encuentra, al menos:
    - Entre la media y dos desviaciones típicas el 75% de las observaciones.
    - Entre la media y tres desviaciones típicas el 89% de las observaciones.
- Esta desigualdad es la famosa ***desigualdad de Tchebychev***.

# Medidas de posicionamiento

- Percentil  $p$ : es el menor valor superior al  $p\%$  de los datos ordenados. Si el número de datos es impar Med es el percentil 50. Es una medida no central usada en estadística que indica el valor de la variable por debajo del cual se encuentra un porcentaje dado de observaciones.
- Cuartiles: aquellos valores que dividen la distribución en cuatro partes iguales.
  - Primer cuartil  $Q_1$ : es el percentil 25.
  - Segundo cuartil  $Q_2$ : es la mediana.
  - Tercer cuartil  $Q_3$ : es el percentil 75.

# Medidas de posicionamiento

- Los percentiles y cuartiles se utilizan para construir medidas de dispersión basadas en datos ordenados:
  - Rango intercuartílico: es la diferencia entre los percentiles 75 y 25.
- Hay muchos tipos **cuantiles**:
  - Percentiles, Quartiles, Deciles, Etc.
- El concepto de mediana se generaliza mediante los cuantiles.
- Un cuantil de **orden  $k$** , será el valor de la variable que deja por debajo de sí una **proporción** de  $k$  observaciones del total  $n$  de todas las observaciones, que se han ordenado previamente en magnitud.
- Es decir habrá  $n \times k$  observaciones con valores menores o igual al cuantil  $k$ .
- Así la media resulta ser el cuantil de orden  $1/2$ .

# Medidas de posicionamiento

- En general los cuantiles son puntos tomados a intervalos regulares de la función de distribución de una variable aleatoria.
- Así cuantil de orden  $p$  de una distribución es el valor de la variable  $x_p$  que indica que por debajo de esta se encuentra un porcentaje dado de observaciones. En este caso estamos en distribuciones, y  $0 < p < 1$ .
  - Ej: el cuantil de orden 0,27 dejaría un 27% de valores están por debajo y el cuantil de orden 0,50 se corresponde con la mediana de la distribución.

# Medidas de posicionamiento

- Resumiendo los tipos de cuantiles más habituales son:
  - Percentiles: dividen a la distribución en cien partes.
  - Cuartiles: dividen a la distribución en cuatro partes (corresponden a los cuantiles 0,25; 0,50 y 0,75).
  - Quintiles: dividen a la distribución en cinco partes (corresponden a los cuantiles 0,20; 0,40; 0,60 y 0,80).
  - Deciles: dividen a la distribución en diez partes.

# Medidas de posicionamiento

- Existen muchos métodos para calcular los percentiles:

Weisstein, Eric W. "Quartile." From MathWorld—A Wolfram Web Resource. <http://mathworld.wolfram.com/Quartile.html>

method	1st quartile	1st quartile	3rd quartile	3rd quartile
	$n$ odd	$n$ even	$n$ odd	$n$ even
Minitab	$\frac{n+1}{4}$	$\frac{n+1}{4}$	$\frac{3n+3}{4}$	$\frac{3n+3}{4}$
Tukey (Hoaglin et al. 1983)	$\frac{n+3}{4}$	$\frac{n+2}{4}$	$\frac{3n+1}{4}$	$\frac{3n+2}{4}$
Moore and McCabe (2002)	$\frac{n+1}{4}$	$\frac{n+2}{4}$	$\frac{3n+3}{4}$	$\frac{3n+2}{4}$
Mendenhall and Sincich (1995)	$\left[ \frac{n+1}{4} \right]$	$\left[ \frac{n+1}{4} \right]$	$\left[ \frac{3n+3}{4} \right]$	$\left[ \frac{3n+3}{4} \right]$
Freund and Perles (1987)	$\frac{n+3}{4}$	$\frac{n+3}{4}$	$\frac{3n+1}{4}$	$\frac{3n+1}{4}$

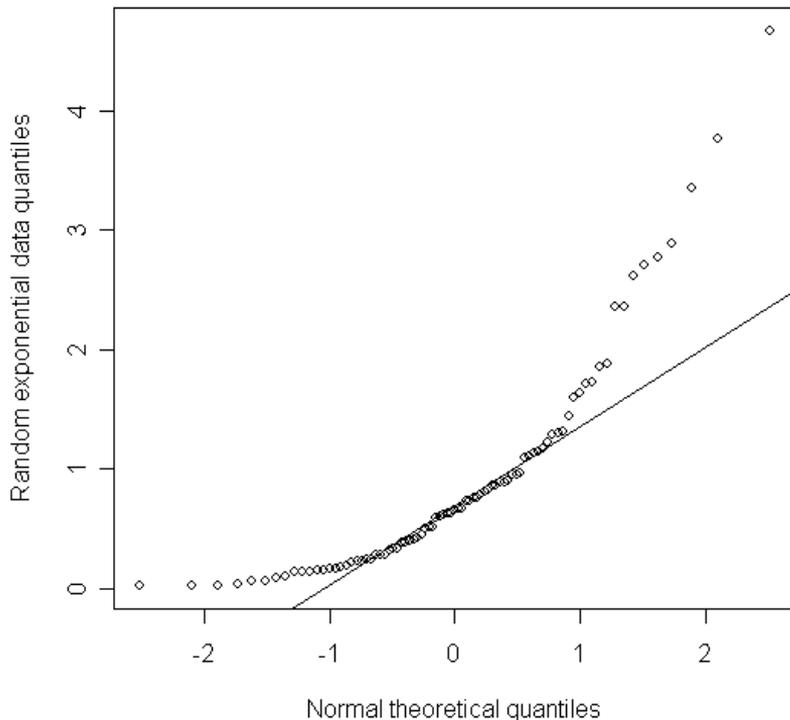
# Visualización datos con medidas de posicionamiento: Gráficos Q-Q

- Los gráficos Q-Q (Q-Q plots) se caracterizan por visualizar de una manera muy rápida y sencilla como se diferencian los datos de dos distribuciones de observaciones.
- Se basan en representar enfrentados en un gráfico x-y los cuantiles de ambas distribuciones. El “Q” viene de cuantil en inglés.
- Si todos los cuantiles son iguales aparecerá la recta  $x=y$  en el gráfico, y significará los dos conjuntos de datos se distribuyen de manera idéntica.
- Generalmente una de la distribuciones es conocida (por ejemplo una normal), para contrastar si los datos observados se ajustan a la distribución conocida.

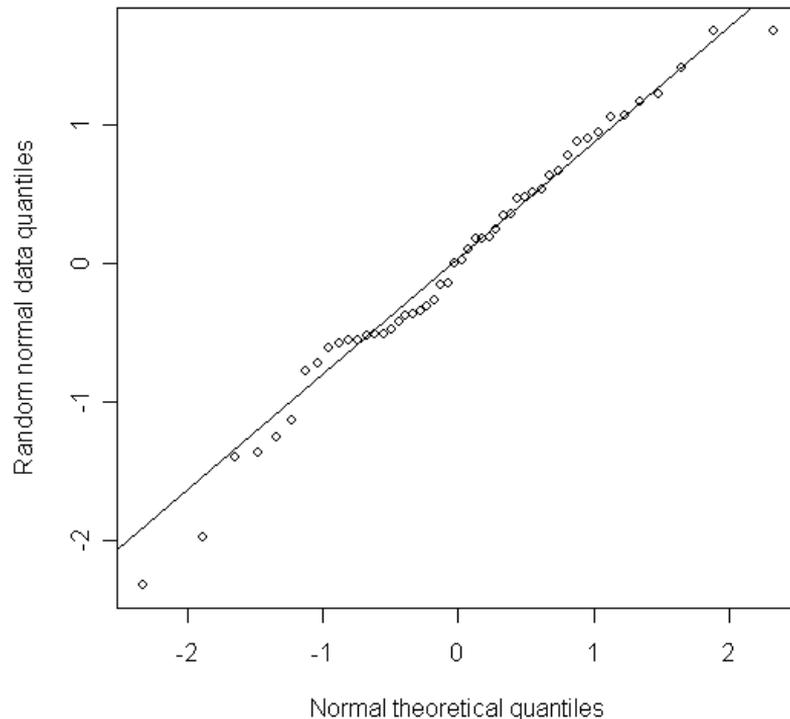
# Visualización datos con medidas de posicionamiento: Gráficos Q-Q

Imágenes extraídas de [https://es.wikipedia.org/wiki/Gr%C3%A1fico\\_Q-Q](https://es.wikipedia.org/wiki/Gr%C3%A1fico_Q-Q)

Normal Q-Q Plot with exponential data



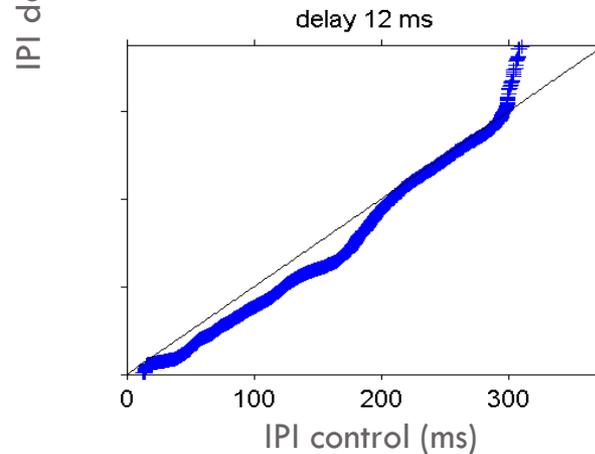
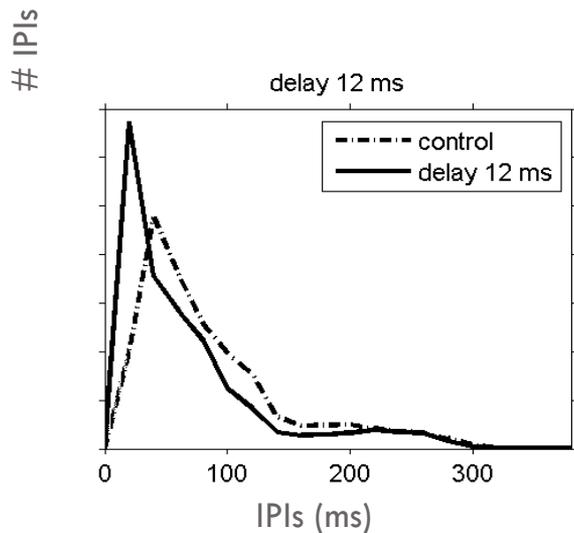
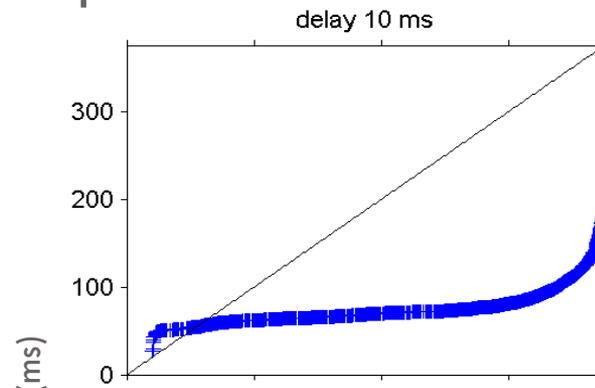
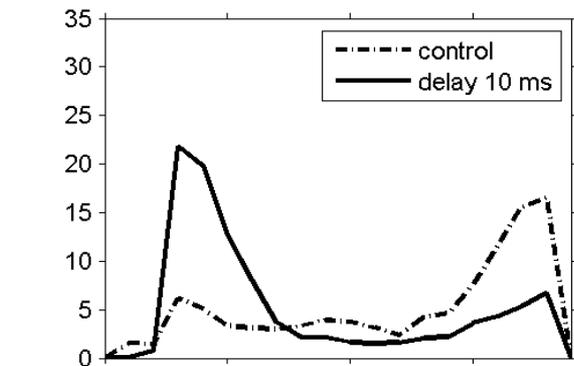
Normal Q-Q Plot



# Visualización datos con medidas de posicionamiento: Gráficos Q-Q

- Para el calculo del cuantil  $k$  y para  $n$  datos se utilizan diversas formulas:
  - La más habitual es  $k/(n+1)$
  - Para graficas simétricas  $(k-a)/(n+1-2a)$ , con  $a \in (0,0.5)$ , con  $a=0$  retomamos la primera, con  $a=0.5$  se suele utilizar para distribuciones normales.
  - $(k-(1/3))/(n+(1/3))$ .
  - $(k-0.3175)/(n+0.365)$ .
  - $(k-0.326)/(n+0.348)$ .
  - Etc.

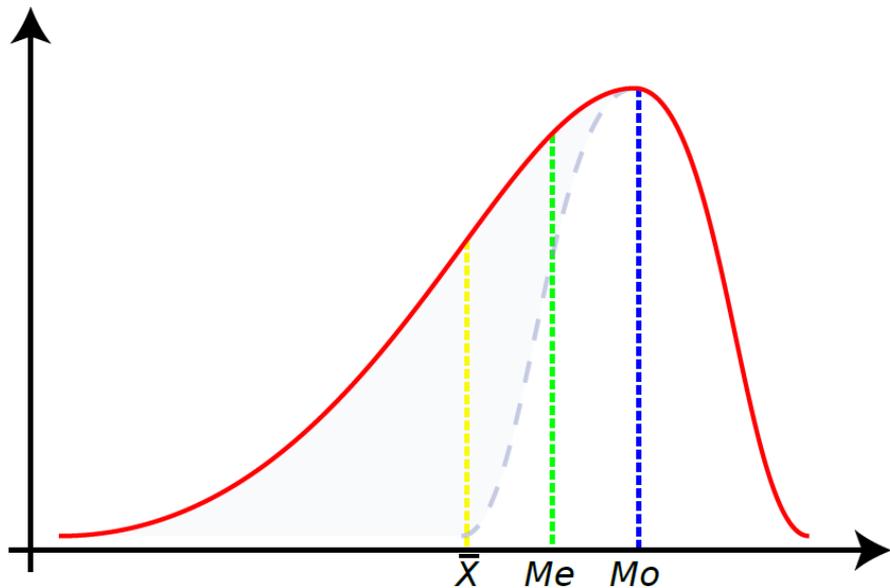
# Visualización datos con medidas de posicionamiento: Gráficos Q-Q



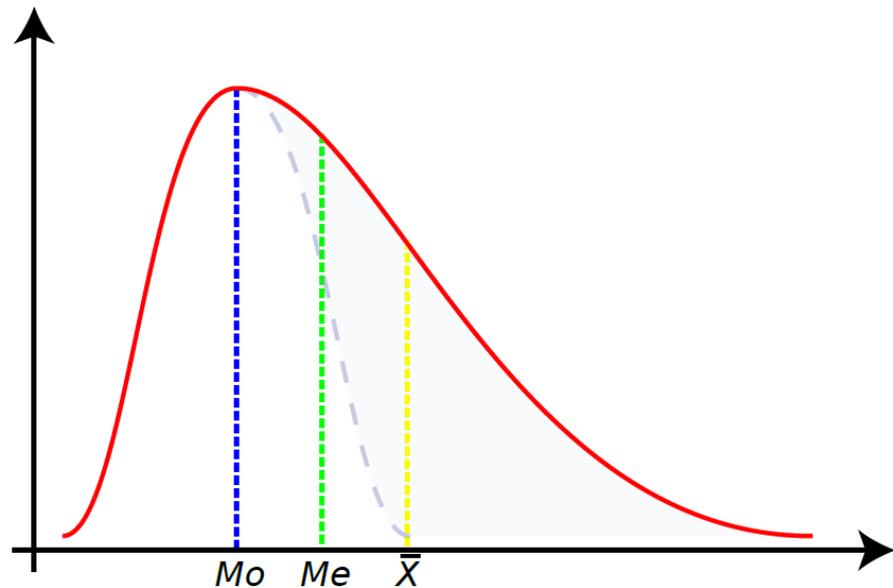
# Medidas de asimetría y curtosis

- Informan de aspectos importantes de la forma de la distribución
- Coeficiente de asimetría (medida adimensional):
  - $CA = \sum (x_i - \langle x \rangle)^3 / ns^3$
  - En un conjunto de datos simétricos respecto a la media el sumatorio es nulo.
  - Si CA es negativo la distribución se alarga para valores inferiores a la media (la cola de la distribución se extiende a la izquierda). La masa de la distribución se concentra en el lado derecho.
  - Si CA es positivo la distribución se alarga para valores superiores a la media (la cola de la distribución se extiende a la derecha). La masa de la distribución se concentra en el lado izquierdo.

# Medidas de asimetría y curtosis



Asimetría negativa



Asimetría positiva

Imagen extraída de [https://es.wikipedia.org/wiki/Asimetr%C3%ADa\\_estad%C3%ADstica](https://es.wikipedia.org/wiki/Asimetr%C3%ADa_estad%C3%ADstica)

# Medidas de asimetría y curtosis

- **Coefficiente de curtosis (apuntamiento y aplastamiento):**
  - $CA_p = \sum(x_i - \langle x \rangle)^4 / ns^4$ , se suele definir como  $CA_p = \sum(x_i - \langle x \rangle)^4 / ns^4 - 3$  ya que la curtosis de la normal es 3, y así la curtosis se expresa en función de la normal.
  - Mide como la frecuencia relativa de unos datos observados se reparte entre el centro y los extremos. A la curtosis también se le llama *apuntamiento*.
  - De otra forma se puede medir la heterogeneidad y presencia de valores atípicos de la distribución:
    - Cuando el apuntamiento es muy bajo se dice que la distribución es heterogénea.
    - Sin embargo, cuanto el apuntamiento es muy alto indica la presencia de valores atípicos en la distribución.
    - Valores intermedios de este coeficiente indica la presencia de homogeneidad en la muestra.
  - También se puede ver este coeficiente en comparación con la distribución normal:
    - Si  $CA_p < 0$  la distribución es más aplanada que la normal (platicúrtica).
    - Si  $CA_p = 0$  la distribución es igual de apuntada que la normal (mesocúrtica).
    - Si  $CA_p > 0$  la distribución es más apuntada que la normal (leptocúrtica).

# Medidas de asimetría y curtosis

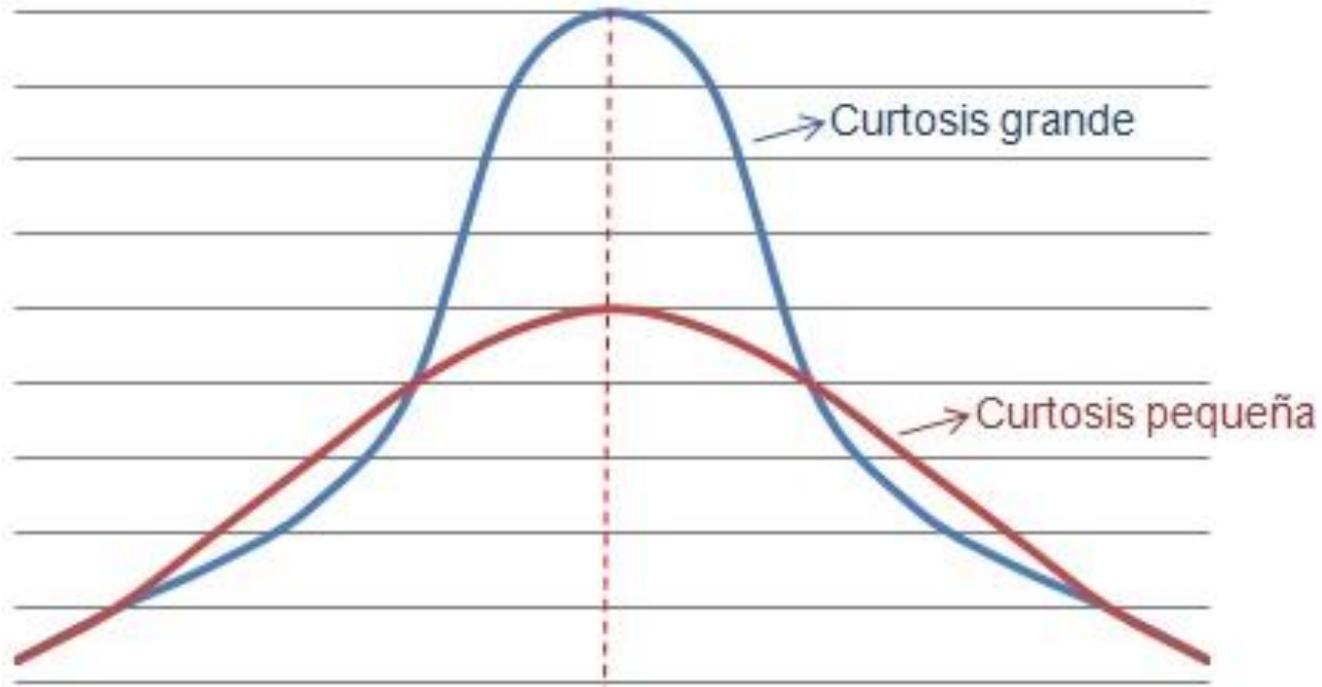


Imagen extraída de <http://www.universoformulas.com/estadistica/descriptiva/curtosis/>

# Datos atípicos y diagramas caja (boxplots)

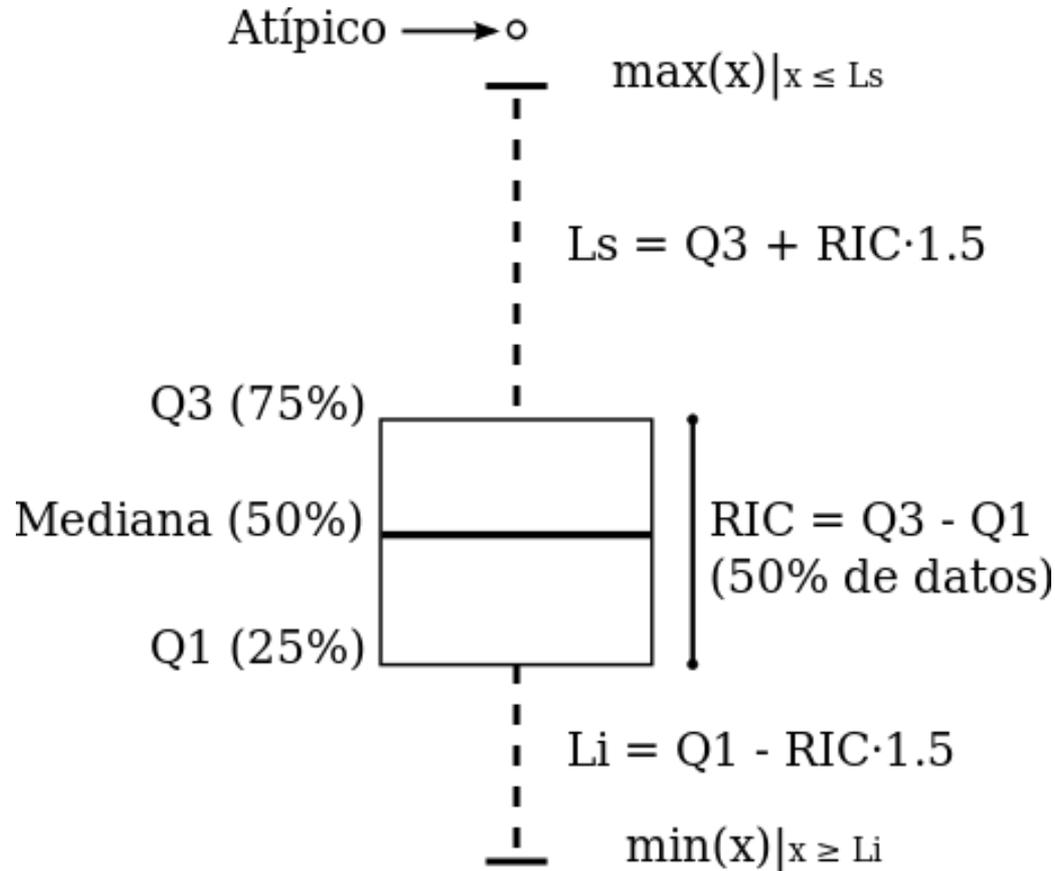
- Datos atípicos: es muy frecuente que los datos presenten cierta observaciones que se ha han medido erróneamente o que se han transcrito mal.
  - Es un histograma pueden ser los datos extremales aislados, pero no solo pueden estar en los extremos.
  - Para detectarlos se utilizan valores de centralización y dispersión que estén poco afectados por estos (Mediana y MEDA).
  - Se suele utilizar esta regla para considerar valores sospechosos:
    - $x > \text{Med} \pm (4.5 \times \text{MEDA})$
  - Este criterio se utiliza mucho pero no tiene en cuenta la asimetría de la distribución. Sin embargo las medidas cuantiles si lo tiene en cuenta y se pueden utilizar par esto.
  - Así considerando los cuartiles y el rango intercuartílico, se suele utilizar esta regla para considerar valores atípicos:
    - $x < Q_1 - 1.5(Q_3 - Q_1)$
    - $x > Q_3 + 1.5(Q_3 - Q_1)$

# Datos atípicos y diagramas caja (boxplots)

- Diagramas caja (boxplots): representación semigráfica de una distribución que muestra sus características principales, señalando los posible datos atípicos (ejercicio 2.5 del libro). El cálculo:
  - Se ordenan los datos y se sacan los cuartiles.
  - Se pintan los cuartiles extremales en un rectángulo con la mediana como separación en el interior.
  - Se calculan los límites admisibles de valores atípicos:
    - $LI < Q_1 - 1.5(Q_3 - Q_1)$
    - $LS > Q_1 + 1.5(Q_3 - Q_1)$
  - Se pintan dos líneas que salen del rectángulo indicando estos límites.

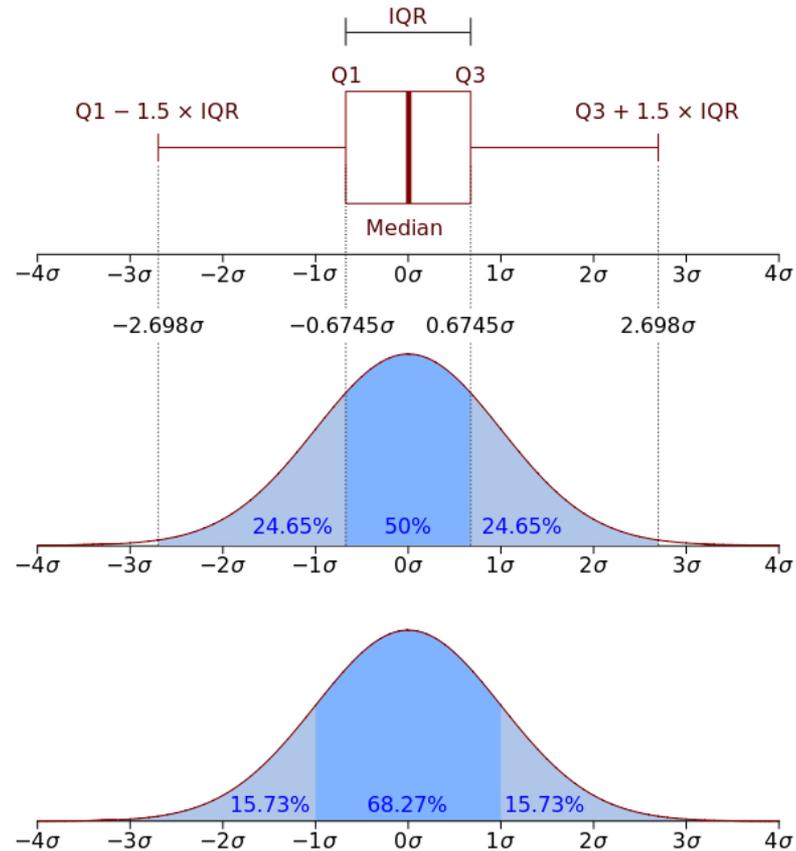
# Datos atípicos y diagramas caja (boxplots)

Imagen extraída de [https://es.wikipedia.org/wiki/Diagrama\\_de\\_caja](https://es.wikipedia.org/wiki/Diagrama_de_caja)



# Datos atípicos y diagramas caja (boxplots)

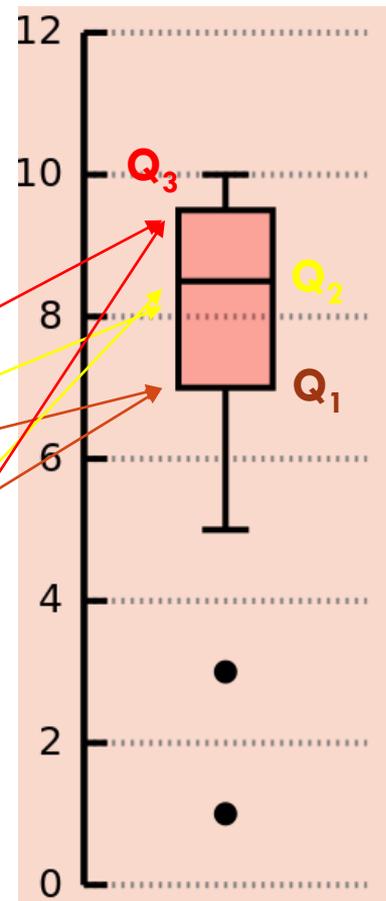
Imagen extraída de [https://en.wikipedia.org/wiki/Box\\_plot](https://en.wikipedia.org/wiki/Box_plot)



# Datos atípicos y diagramas caja (boxplots)

Imagen extraída de <https://de.wikipedia.org/wiki/Boxplot>

$i$	1	2	3	4	5
$x_i$	9	6	7	7	3
	9	10	1	8	7
	9	9	8	10	5
	10	10	9	10	8
$x_i$ (sort)	7	3	5	6	7
	7	7	8	8	8
	9	9	9	9	9
	10	10	10	10	10



# Datos: Descripción conjunta de variables

## Descripción conjunta de varias variables.

- Distribuciones de frecuencia multivariantes
- Medidas de dependencia lineal
- Recta de Regresión
- Matriz de varianzas

## Distribuciones de frecuencia multivariantes

- Uno de los objetivos del análisis estadístico es encontrar las relaciones que existen entre un grupo de variables.
- Ahora supondremos que nuestro conjunto de datos contiene los valores de las variables  $(x, y)$ , que se han medido conjuntamente en una población.
- El análisis que veremos aquí se generaliza por cualquier número de variables.

# Distribuciones de frecuencia multivariantes: D. Conjunta

- La distribución conjunta de frecuencias de dos variables ( $x, y$ ) es una tabla que representa los valores observados de ambas variables, y las frecuencias relativas de aparición de cada pareja de valores.

- La suma de las frecuencias relativas

- $\sum_i \sum_j f_r(x_i, y_j) = 1$

- Frecuencias absoluta de una pareja:

- $f_r(x_i, y_j) \times n$ , donde  $n=20$ .

- $n_1(70,175) = 0.05 \times 20 = 1$

- $n_5(70,165) = 0.2 \times 20 = 4$

Pesos $X_i$	Altura $Y_j$	$n_i(x_i, y_j)$	$f_r(x_i, y_j)$
70	175	1	0.05
65	160	3	0.15
85	180	2	0.1
60	155	3	0.15
70	165	4	0.2
75	180	2	0.1
90	185	1	0,05
80	175	1	0.05
60	160	2	0.1
70	170	1	0.05
		$n=20$	1

# Distribuciones de frecuencia multivariantes:

## Distribución Marginal

- La distribución marginal de frecuencias de una variable se obtiene al estudiar esta aislada, e independiente del resto, así:
  - La distribución marginal de  $\mathbf{x}$  se obtiene como  $f_r(\mathbf{x}_i) = \sum_j f_r(\mathbf{x}_i, \mathbf{y}_j)$ .
  - Y análogamente la de  $\mathbf{y}$  como  $f_r(\mathbf{y}_j) = \sum_i f_r(\mathbf{x}_i, \mathbf{y}_j)$ .

# Distribuciones de frecuencia multivariantes: D. Condicional

- **Distribución condicional:** la distribución de  $y$  para  $x=x_i$  es la distribución **univariante** de la variable  $y$  que se obtiene considerando sólo los elementos que tienen para la variable  $x$  el valor  $x_i$ :
  - $f_r(y_i | x_i) = f_r(x_i, y_i) / f_r(x_i)$ .
  - Se normaliza por  $f_r(x_i)$ : de esta forma se garantiza que la suma de frecuencias relativas para todos los valores de la variable  $y$  es 1:
    - $\sum_i f_r(y_i | x_i) = \sum_i f_r(x_i, y_i) / f_r(x_i) = f_r(x_i) / f_r(x_i) = 1$ .
- La distribución condicional de  $y$  para  $x=x_i$  se interpreta como la distribución de la característica  $y$  en los elementos de la población que tienen como característica  $x$  el valor de  $x_i$ .

# Distribuciones de frecuencia multivariantes: D. Condicional

- La diferencia con la distribución marginal es clara: La distribución marginal en  $\mathbf{y}$  tiene en cuenta la distribución de  $\mathbf{y}$  en todos los elementos con independencia del valor que en ellos tenga la característica  $\mathbf{x}$ . En cambio la condicional fija un conjunto de valores de  $\mathbf{x}=\mathbf{x}_j$ .
- Se puede deducir la distribución marginal de la condicional:
  - $f_r(\mathbf{y}) = \sum_i f_r(\mathbf{y} | \mathbf{x}_i) \times f_r(\mathbf{x}_i)$ , i.e. la frecuencia de  $\mathbf{y}$  en la población total se obtiene ponderando su frecuencia en las subpoblaciones definidas por los distintos valores de  $\mathbf{x}$ .
- De igual forma la distribución conjunta se obtiene de la condicional, si conocemos todas las distribuciones marginales y condicionales para cada variable:
  - $f_r(\mathbf{x}_i, \mathbf{y}_i) = f_r(\mathbf{y}_i | \mathbf{x}_i) \times f_r(\mathbf{x}_i)$ .

# Distribuciones de frecuencia multivariantes:

## Representaciones Gráficas

- La representación de las observaciones de dos o tres variables enfrentadas es lo que se llaman diagramas de dispersión.
- En el se puede observar muy rápidamente la posible existencia de las variables:
  - Relación lineal positiva.
  - Relación lineal negativa.
  - Falta de relación.
  - Otros tipos de relación.
- Cuando tenemos varias variables se pueden replantear como una matriz de figuras de diagramas de dispersión y en la diagonal los histogramas de las variables.

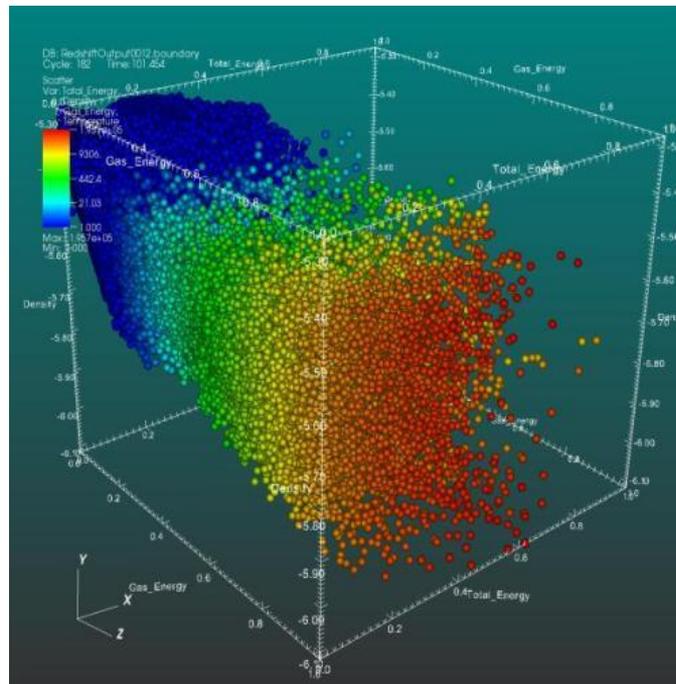
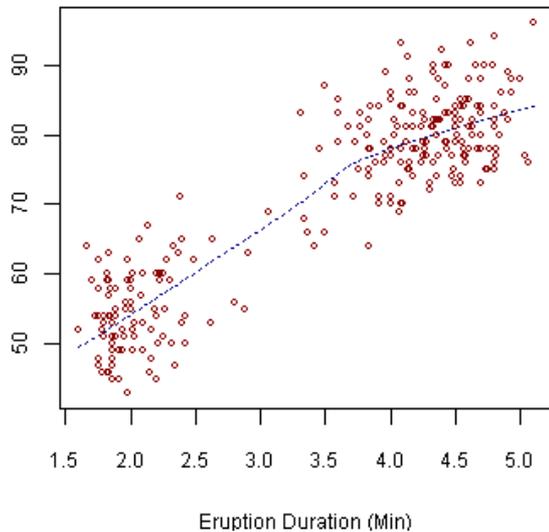
# Distribuciones de frecuencia multivariantes: Representaciones Gráficas

Imágenes extraídas de

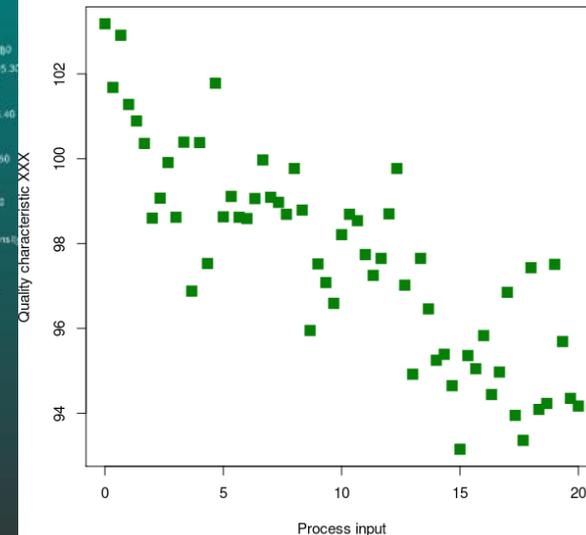
[https://en.wikipedia.org/wiki/Scatter\\_plot](https://en.wikipedia.org/wiki/Scatter_plot)

### Old Faithful Eruptions

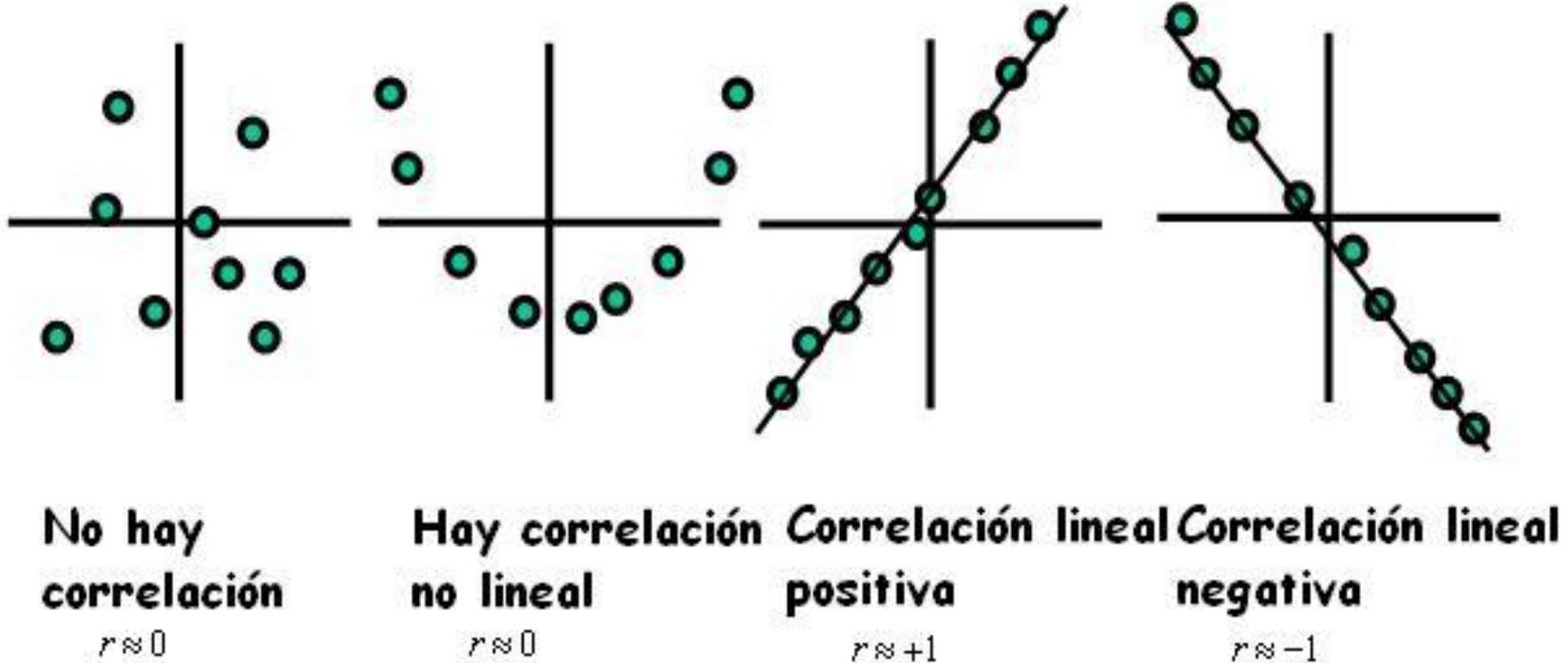
Waiting Time Between Eruptions (Min)



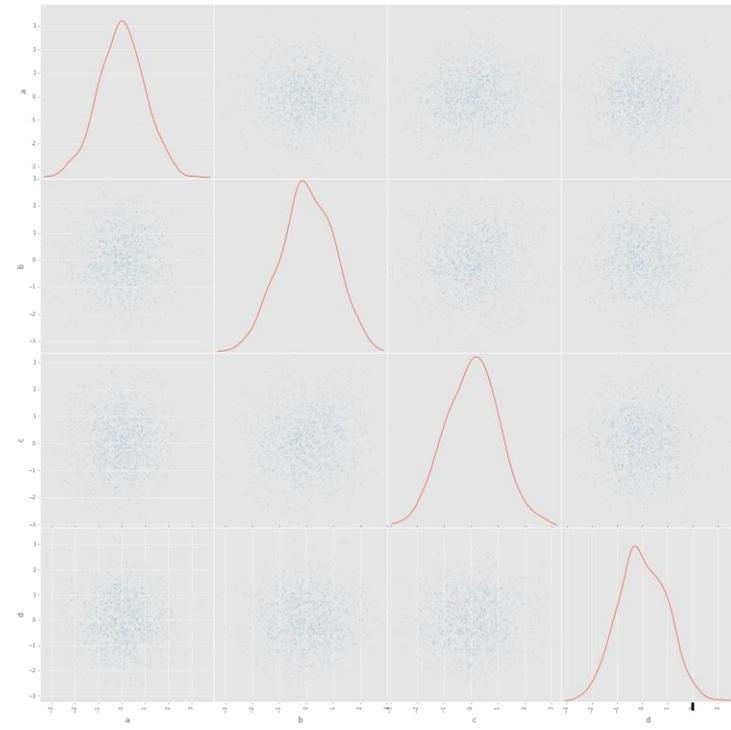
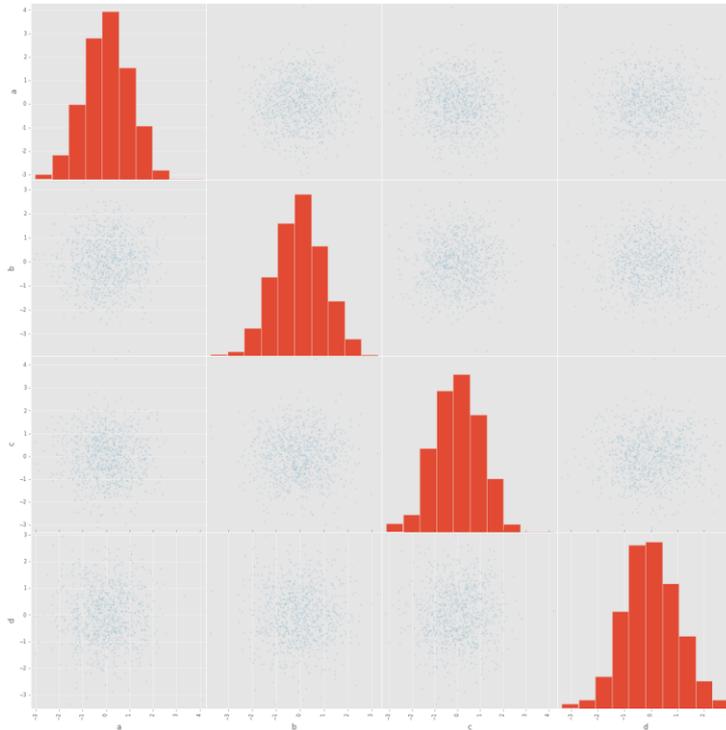
### Scatterplot for quality characteristic XXX



# Distribuciones de frecuencia multivariantes: Representaciones Gráficas



# Distribuciones de frecuencia multivariantes: Representaciones Gráficas



# Medidas de dependencia lineal: Covariancia

- La covariancia es la medida más utilizada como una medida descriptiva de la posible relación lineal que puede existir entre un par de variables:
  - $\text{Cov}(x, y) = \sum_i ((x_i - \langle x \rangle)(y_i - \langle y \rangle)) / n$ , siendo este sumatorio sobre todas las posibles  $n$  parejas de valores  $(x, y)$ .
  - Se puede deducir que una expresión equivalente es:
    - $\text{Cov}(x, y) = \{\sum_i x_i y_i / n\} - \langle x \rangle \langle y \rangle$ .
  - Si los datos están agrupados por clases, con frecuencia relativa de cada clase  $f_r(\mathbf{x}_i, \mathbf{y}_i)$ , entonces:
    - $\text{Cov}(x, y) = \sum_i (x_i - \langle x \rangle)(y_i - \langle y \rangle) f_r(\mathbf{x}_i, \mathbf{y}_i)$ .

## Medidas de dependencia lineal: Covariancia

- La covarianza fue introducida por K. Pearson para medir la relación lineal entre las variables, como en la figura anterior.
- El **signo positivo** de la covarianza indica que cuando una variable esta por encima de su media, es esperable que la otra variable este por encima de la suya también (producto positivo en el sumatorio).
- El **signo negativo** de la covarianza indica que cuando una variable esta por encima de su media, es esperable que la otra variable este por debajo de la suya (producto negativo en el sumatorio).

# Medidas de dependencia lineal: Correlación

- La covarianza depende de las unidades de las variables a analizar:
  - p. ej. si calculamos la covarianza entre la estatura medida en centímetros y el peso en gramos y nos sale de valor  $cov$ , esta misma covarianza aparecería dividida por  $10^5$ , si las medidas utilizadas en la variables son metros y kilogramos.
- La correlación es una medida adimensional que mide la relación lineal entre dos variables, coeficiente de correlación  $r$ :
  - $r = Cov(x, y) / s_x s_y$ , con  $s_x$  y  $s_y$  las desviaciones típicas de las variables  $x$  e  $y$ .
  - Como hemos dividido por un termino que tiene las mismas dimensiones que la  $Cov$ , se hace adimensional.

# Medidas de dependencia lineal: Correlación

- La elección de las desviaciones típicas como factor dividiendo dota a  $r$  de una serie de propiedades interesantes:
  - $r$  tiene el mismo signo que  $\text{Cov}$ .
  - $r$  es adimensional:  $r$  no varia si multiplicamos a  $\mathbf{x}$  por  $k_1$  y a  $\mathbf{y}$  por  $k_2$ , siempre que  $k_1$  y  $k_2$  sean no nulos y del mismo signo.
  - Si existe una relación pura lineal entre  $\mathbf{x}$  e  $\mathbf{y}$ , i.e.  $y=a+bx$ , entonces:
    - $r=1$  sii  $b>0$  o  $r=-1$  sii  $b<0$ .
  - Si no existe un relación lineal pura entonces  $r$  tiene un valor comprendido ente 1 y -1.
- Así el coeficiente de correlación puede resumir un diagrama de dispersión, como hemos visto en la figura anterior.

# Recta de regresión

- Cuando existe una relación lineal entre variables, los puntos se agrupan en una línea en un diagrama de dispersión y la forma natural de describir el sistema es mediante la mejor recta que ajusta el diagrama de dispersión.
- De la misma forma que describimos una variable con su media y dispersión, podemos describir dos variables con una recta, y su dispersión: la que existe entre los puntos y la recta.
- Supongamos que queremos minimizar los errores de la variable  $y$  cuando conocemos el valor de  $x$ , entonces la recta será de la forma:
  - $h(x)=a+bx$

# Recta de regresión

- Así si minimizamos  $\sum (y_i - h(\mathbf{x}))^2 = \sum (y_i - a - b\mathbf{x}_i)^2$ , y obtenemos la recta que mejor se ajusta al diagrama de dispersión lineal.
- Minimizar significa derivar respecto los parámetros a optimizar, i.e. respecto a y b, e igualar a 0 la derivada.
- Por tanto derivamos respecto a ambos coeficientes e igualando a 0, y resolvemos el sistema de 2 ecuaciones con dos incógnitas:
  - $2\sum (y_i - a - b\mathbf{x}_i) (-1) = 0$ ,
  - $2\sum (y_i - a - b\mathbf{x}_i) (-\mathbf{x}_i) = 0$

# Recta de regresión

- Dividiendo por  $n$ , i.e. el número de parejas  $(x_i, y_i)$ , obtenemos:
  - $\langle y \rangle = a + b \langle x \rangle$ , que indica que la recta debe pasar por el centro de la nube,
  - $\sum (x_i y_i) / n = a \langle x \rangle + b \sum x_i^2 / n$ .
- Si eliminamos el término  $a$  de la segunda ecuación restando la primera multiplicada por  $\langle x \rangle$  obtenemos:
  - $\sum_i x_i y_i / n - \langle x \rangle \langle y \rangle = b (\sum (x_i^2 / n) - \langle x \rangle^2)$
- Así  $b = \text{Cov}(x, y) / s_x^2$ , indicando que la pendiente de la recta es la covarianza estandarizada para que tenga unidades de  $x/y$  como corresponde a la pendiente de una recta.

$$s = [\sum (x_i^2 / n) - \langle x \rangle^2]^{1/2}$$
$$\text{Cov}(x, y) = \{ \sum_i x_i y_i / n \} - \langle x \rangle \langle y \rangle$$

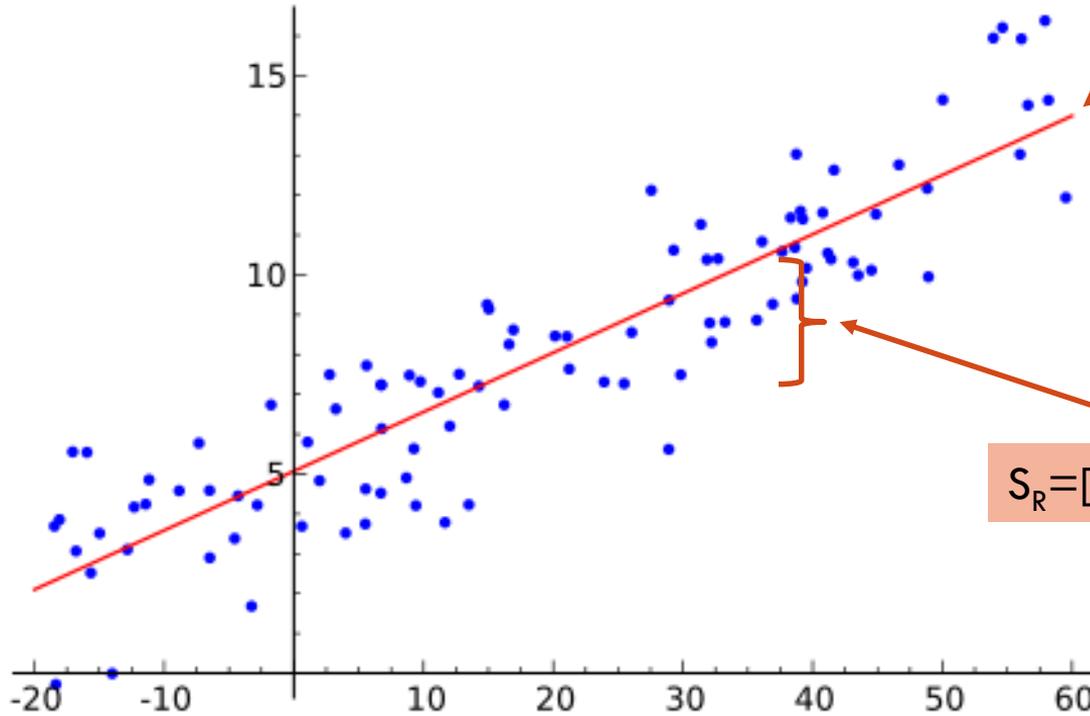
# Recta de regresión

- Sustituyendo en la ecuación de la recta  $a = \langle y \rangle - b \langle x \rangle$  y  $b = \text{Cov}(x, y) / s_x^2$ , se obtiene:
  - $h(x) = \langle y \rangle + (\text{Cov}(x, y) / s_x^2)(x - \langle x \rangle)$ , i.e. la regresión lineal (su nombre por F. Galton y su famosa regresión de la media).
- La medida de variabilidad de los datos respecto a la recta calculada (al igual que calculábamos la varianza en la descripción de un sola variable) se denomina desviación típica residual:
  - $S_R = [ \sum (y_i - h(x_i))^2 / n ]^{1/2}$ , es el promedio de las desviaciones verticales de la recta con los datos.



# Recta de regresión

$$h(x) = \langle y \rangle + \text{Cov}(x, y) / s_x^2 (x - \langle x \rangle)$$



$$S_R = \left[ \sum (y_i - h(x_i))^2 / n \right]^{1/2}$$

Parte de la foto sacada de  
[https://en.wikipedia.org/wiki/Linear\\_regression](https://en.wikipedia.org/wiki/Linear_regression)

# Vector de medias y matriz de varianzas y covarianzas

- El vector de medias por ejemplo para una variable tridimensional:

$$\text{➤ } \langle X \rangle = \begin{bmatrix} \langle x \rangle \\ \langle y \rangle \\ \langle z \rangle \end{bmatrix} = \frac{1}{n} \begin{bmatrix} \sum x_i \\ \sum y_i \\ \sum z_i \end{bmatrix} = \frac{1}{n} \sum X_i$$

$$\text{➤ } \text{Y en general } \langle X \rangle = 1/n \sum X_i$$

- La matriz de covarianzas para dos variables:

$$\text{➤ } M = \begin{bmatrix} S_x^2 & \text{cov}(x, y) \\ \text{cov}(y, x) & S_y^2 \end{bmatrix}, \text{ es simétrica ya que } \text{cov}(x, y) = \text{cov}(y, x)$$

# Vector de medias y matriz de varianzas y covarianzas

- Para caso k-dimensional variable, si definimos  $s_i^2$  como la varianza de la variable i y  $s_{ij}$  como la covarianza de la variable i con la j, la matriz de covarianzas quedaría:

- $$M = \begin{bmatrix} S_1^2 & \cdots & S_{1k} \\ \vdots & \ddots & \vdots \\ S_{k1} & \cdots & S_k^2 \end{bmatrix}$$

- La matriz de covarianzas se calcula conociendo el vector de medias como:

- $$M = \frac{1}{n} \sum (X_i - \langle X \rangle)(X_j - \langle X \rangle),$$
 por ejemplo para tres variables:

- $$M = \frac{1}{n} \sum \begin{bmatrix} x_i - \langle x \rangle \\ y_i - \langle y \rangle \\ z_i - \langle z \rangle \end{bmatrix} [x_i - \langle x \rangle \quad y_i - \langle y \rangle \quad z_i - \langle z \rangle],$$
 i.e.

- $$M = \begin{bmatrix} (x_i - \langle x \rangle)^2 & (x_i - \langle x \rangle)(y_i - \langle y \rangle) & (x_i - \langle x \rangle)(z_i - \langle z \rangle) \\ (y_i - \langle y \rangle)(x_i - \langle x \rangle) & (y_i - \langle y \rangle)^2 & (y_i - \langle y \rangle)(z_i - \langle z \rangle) \\ (z_i - \langle z \rangle)(x_i - \langle x \rangle) & (z_i - \langle z \rangle)(y_i - \langle y \rangle) & (z_i - \langle z \rangle)^2 \end{bmatrix}$$