

Tema 1 Análisis exploratorio de datos

José R. Berrendero

Departamento de Matemáticas
Universidad Autónoma de Madrid

José Ramón Berrendero Díaz

Correo electrónico: `joser.berrendero@uam.es`

Teléfono: 91 497 66 90

Despacho: Módulo 08 - Despacho 210

Página web: <http://www.uam.es/joser.berrendero>

Ejemplo de introducción: contaminación por mercurio en el pescado

- ▶ El agua de los ríos contiene pequeñas concentraciones de mercurio que se pueden ir acumulando en los tejidos de los peces.
- ▶ Se ha realizado un estudio en los ríos Wacamaw y Lumber en Carolina del Norte (EE.UU.), analizando la cantidad de mercurio que contenían 171 ejemplares capturados de una cierta especie de peces.
- ▶ Los datos obtenidos se encuentran en el fichero `mercurio.txt` (formato texto) o en el fichero `mercurio.sav` (formato SPSS).

Variables

Nombre variable	Descripción
RIO	Código del río (0=Lumber, 1=Wacamaw)
ESTACION	Código de la estación (de 0 a 16)
LONG	Longitud (en cm) del pez
PESO	Peso (en g) del pez
CONC	Concentración (en ppm) de mercurio

mercuro.sav [Conjunto de datos1] - Editor de datos SPSS

Archivo Edición Ver Datos Transformar Analizar Gráficos Utilidades Ventana ?

74 :

	RIO	ESTACION	LONG	PESO	CONC	VAR	VAR	VAR	VAR	VAR	VAR	VAR	VAR	VAR	VAR	VAR
46	.00	3.00	40.50	863.00	1.10											
47	.00	3.00	38.50	773.00	1.30											
48	.00	3.00	40.00	898.00	1.20											
49	.00	3.00	44.50	1301.00	1.60											
50	.00	3.00	50.50	2021.00	1.20											
51	.00	3.00	50.00	1883.00	1.20											
52	.00	3.00	47.00	1546.00	1.20											
53	.00	4.00	26.00	234.00	.28											
54	.00	4.00	30.00	323.00	.51											
55	.00	4.00	30.00	430.00	.27											
56	.00	4.00	29.00	353.00	.35											
57	.00	4.00	31.00	402.00	.57											
58	.00	4.00	41.00	1042.00	.66											
59	.00	4.00	36.50	723.00	.49											
60	.00	4.00	38.00	709.00	.51											
61	.00	4.00	44.00	1455.00	.73											
62	.00	4.00	27.50	308.00	.94											
63	.00	4.00	39.00	879.00	1.40											
64	.00	4.00	46.00	1396.00	2.70											
65	.00	4.00	51.50	2389.00	1.80											
66	.00	4.00	56.00	3421.00	3.10											
67	.00	5.00	30.50	400.00	.50											
68	.00	5.00	42.00	1049.00	1.20											
69	.00	5.00	34.00	678.00	.55											
70	.00	6.00	50.50	2228.00	3.50											
71	.00	6.00	44.50	1853.00	2.70											
72	.00	6.00	30.00	380.00	1.40											
73	.00	6.00	35.00	620.00	.73											
74	1.00	7.00	30.40	320.00	.51											
75	1.00	7.00	25.20	217.00	.23											
76	1.00	7.00	36.00	574.00	1.60											
77	1.00	7.00	29.90	298.00	.48											
78	1.00	7.00	34.70	491.00	.58											
79	1.00	7.00	34.70	491.00	.58											
80	1.00	7.00	31.70	407.00	.50											
81	1.00	7.00	34.50	496.00	.58											
82	1.00	7.00	40.10	805.00	1.40											
83	1.00	7.00	30.70	315.00	.44											
84	1.00	8.00	33.00	518.00	.55											
85	1.00	8.00	40.50	970.00	.84											
86	1.00	8.00	32.00	421.00	.65											
87	1.00	8.00	37.00	673.00	.76											
88	1.00	8.00	35.50	550.00	.65											

Vista de datos Vista de variables

SPSS El procesador está preparado

Problemas de interés relacionados con estos datos

- ▶ Resumir la información que contienen con unas pocas cifras o gráficos.
- ▶ ¿Que valores toma cada variable? ¿Cuáles son los más frecuentes? ¿Hay grandes diferencias entre ellos?
- ▶ ¿Es significativamente más alta la concentración de mercurio en un río que en otro?
- ▶ ¿Existe relación entre la concentración de mercurio y la longitud o el peso del pez?
- ▶ ¿Depende la concentración de mercurio de la estación en la que ha sido capturado el pez?

Temario

- ▶ Análisis exploratorio de datos
 - ▶ Introducción al programa SPSS
- ▶ Nociones elementales de inferencia estadística.
 - ▶ La distribución normal
- ▶ Contraste de hipótesis
- ▶ Regresión lineal simple
- ▶ Análisis de la varianza (en función del tiempo disponible)

Bibliografía

- ▶ Freedman, D., Pisani, R., Purves, R. y Adhikari, A. (1993). Estadística. Antoni Bosch ed., Barcelona.
- ▶ de la Horra, J. (2003). Estadística Aplicada. Ediciones Díaz de Santos, Madrid.
- ▶ Milton, J.S. (2001). Estadística para Biología y Ciencias de la Salud. Mc- Graw Hill Interamericana, Madrid.
- ▶ Moore, D.S. (1998). Estadística aplicada básica. Antoni Bosch ed., Barcelona.
- ▶ Rosner, B. (2006). Fundamentals of Biostatistics. Thomson Brooks/Cole.

Estructura del Tema 1

- ▶ Tipos de variables
- ▶ Distribución de una variable
- ▶ Representación gráfica de la distribución
- ▶ Medidas numéricas para resumir la distribución
- ▶ Correlación
- ▶ Transformaciones: estandarización y transformación logarítmica

Tipos de variables

1. **Variables cualitativas:** Describen cualidades o atributos (ej. color del pelo).
2. **Variables cuantitativas discretas:** Toman un número pequeño de valores, normalmente enteros (ej. número de hijos).
3. **Variables cuantitativas continuas:** Toman valores en un intervalo (ej. tiempo hasta que llega un autobús).

En los datos sobre contenido de mercurio, ¿de qué tipo es cada una de las variables?

En general, la técnica estadística adecuada para analizar una variable depende de su tipo.

Introducción

La *estadística* tiene por objetivo extraer conocimiento a partir de información (principalmente) numérica.

El *análisis exploratorio de datos* (o estadística descriptiva) tiene por objetivo identificar las principales características de un conjunto de datos mediante un número reducido de gráficos y/o números.

Los conjuntos de datos que vamos a considerar proceden de medir una o más *variables* en un conjunto de *individuos*.

Para describir un conjunto de datos se comienza con un análisis individual de cada variable y posteriormente se estudian las relaciones entre variables.

Se suele comenzar con representaciones gráficas y posteriormente se calculan resúmenes numéricos.

Distribución de una variable

La *distribución de una variable* viene determinada por los valores que toma esa variable y la frecuencia con la que los toma.

La *frecuencia absoluta* de un valor (o de un intervalo) es el número de individuos para los que la variable toma ese valor (o pertenece a ese intervalo).

La *frecuencia relativa* es igual a la frecuencia absoluta dividida por el número de datos n .

La frecuencia relativa siempre es un número entre 0 y 1.

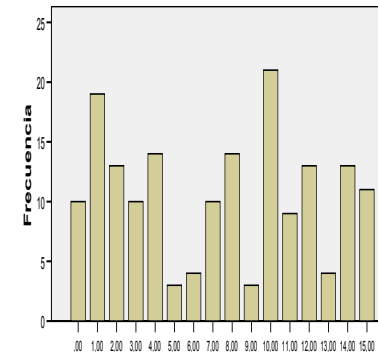
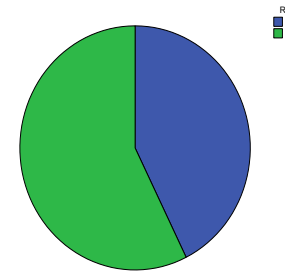
Aspectos interesantes de una distribución

- ▶ Su *posición*: en torno a qué valor central toma valores la variable.
- ▶ Su *dispersión*: el grado de concentración de los valores que toma la variable alrededor de su posición central.
- ▶ Su *forma*: por ejemplo, la simetría, es decir, si los valores se reparten de la misma forma a uno y otro lado del centro.

Piensa en dos conjuntos de 5 datos que tengan:

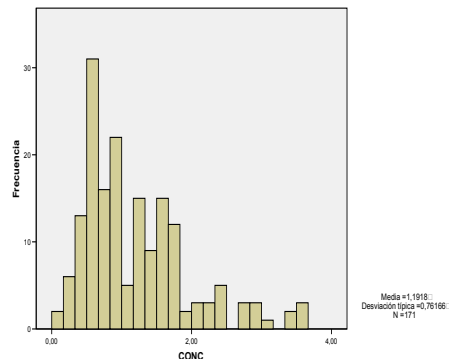
- (a) La misma posición y distinta dispersión.
- (b) La misma dispersión y distinta posición.

Sectores o barras (sólo datos cualitativos o discretos)



Histogramas

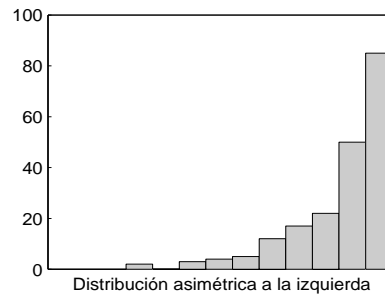
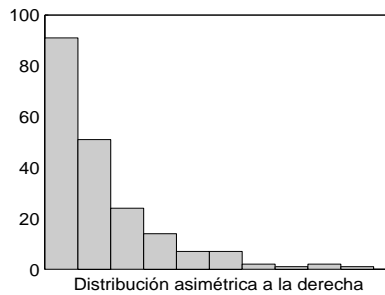
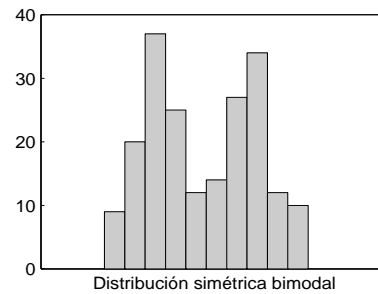
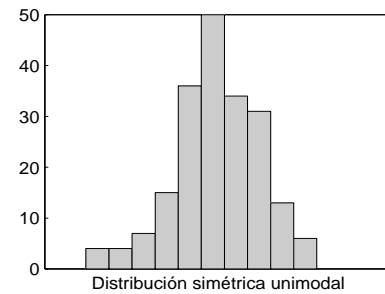
- ▶ Se divide el rango de los datos en un número adecuado de intervalos.
- ▶ Sobre cada intervalo se dibuja un rectángulo cuya área es proporcional a la frecuencia (relativa o absoluta) de datos en el intervalo.



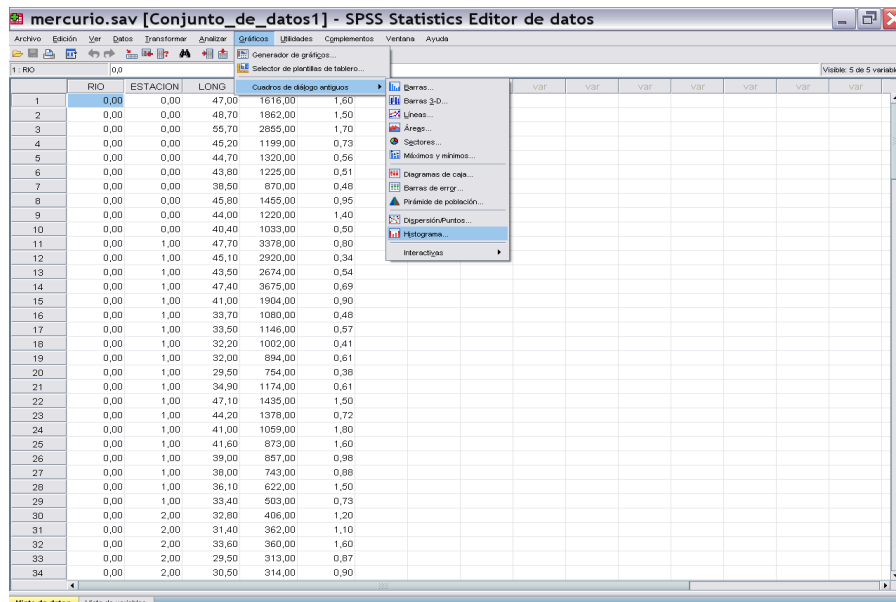
Aspectos a tener en cuenta para interpretar un histograma

- ▶ Normalmente la base de todos los rectángulos es la misma por lo que la altura es proporcional a la frecuencia.
- ▶ Identificar si se han usado frecuencias absolutas o relativas.
- ▶ ¿Cuántas modas hay?
- ▶ ¿Hay algún dato atípico en relación al resto?
- ▶ ¿Es simétrica la distribución?
- ▶ En caso de asimetría, ¿es asimétrica a la izquierda o a la derecha?
- ▶ ¿En torno a qué valor aproximado están centrados los datos?
- ▶ ¿Están muy dispersos los datos en torno a este centro?

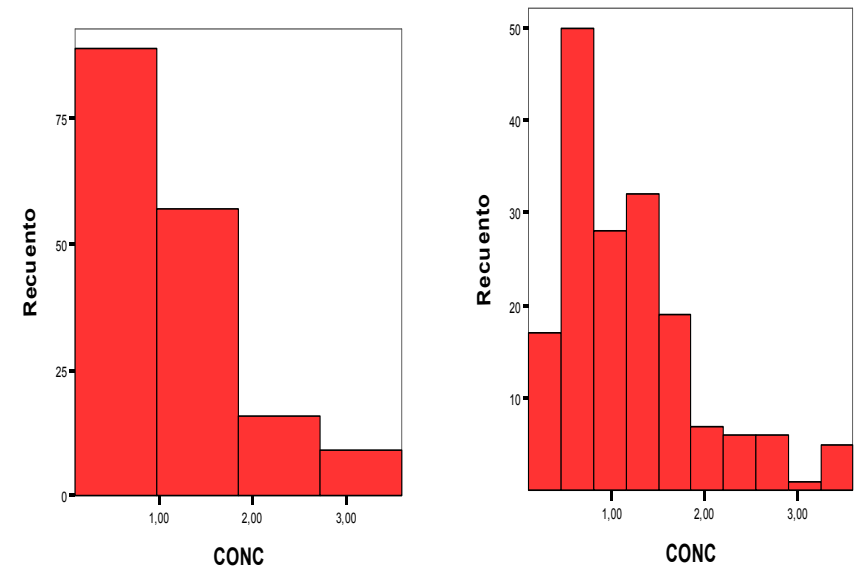
Tipos de simetría



Con SPSS



La forma depende del número de intervalos



Medidas numéricas de posición: la media aritmética

La medida de posición más conocida es la **media aritmética o promedio** de los datos:

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

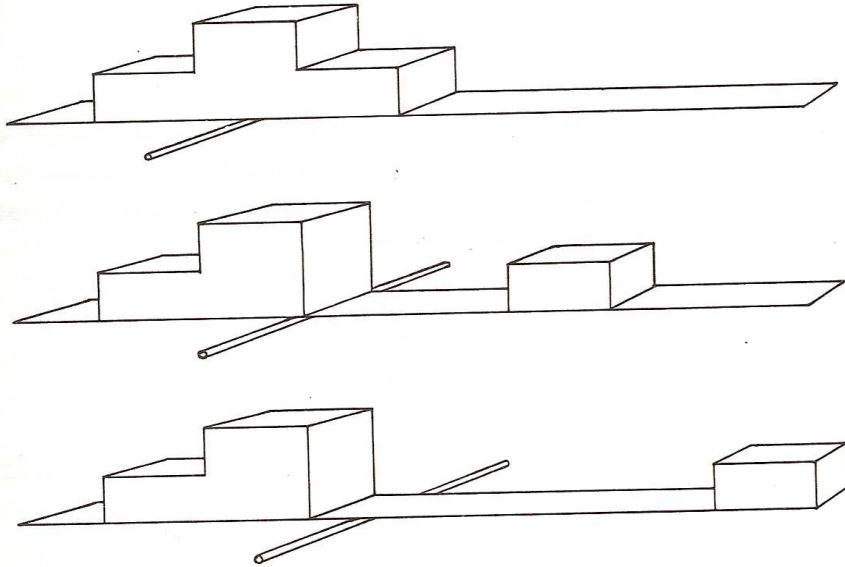
Algunas propiedades:

- ▶ La suma de las desviaciones a la media siempre es igual a cero:

$$(x_1 - \bar{x}) + (x_2 - \bar{x}) + \dots + (x_n - \bar{x}) = 0.$$

- ▶ Si la distribución es muy asimétrica, la media puede distorsionar nuestra percepción de cómo son los datos.
- ▶ La media es muy sensible a la existencia de datos atípicos en los datos.

Posición de la media en un histograma



Medidas numéricas de posición: la mediana

Una medida alternativa de posición es la **mediana**. Para calcular la mediana:

- ▶ Se ordenan los datos de menor a mayor.
- ▶ Si el número de datos es impar, la mediana es el dato que ocupa la posición central.
- ▶ Si el número de datos es par, la mediana es la media de los dos datos centrales.

La mediana es *más robusta* que la media pero hace un uso menos eficiente de la información contenida en los datos.

Relación entre la simetría de una distribución y la posición relativa entre la media y la mediana.

Medidas de dispersión: el rango y los cuartiles

Una medida de dispersión muy sencilla es el *rango o recorrido* de los datos: el valor máximo menos el mínimo.

El rango sólo depende de los datos extremos por lo que no es muy conveniente.

Mejores propiedades tienen los cuartiles y el rango intercuartílico:

- ▶ El **primer cuartil**, Q_1 , es la mediana de los datos menores que la mediana.
- ▶ El **tercer cuartil**, Q_3 , es la mediana de los datos mayores que la mediana.
- ▶ El **rango, recorrido o amplitud intercuartílica** es la diferencia entre los dos cuartiles anteriores: $Q_3 - Q_1$.

De acuerdo con las anteriores definiciones, responde a las siguientes cuestiones:

¿Qué porcentaje de datos hay...

- (a) ... entre Q_1 y Q_3 ?
- (b) ... a la izquierda de Q_1 ?
- (c) ... a la derecha de Q_3 ?
- (d) ... entre el mínimo y Q_3 ?

Una descripción útil de un conjunto de datos viene dada por los cinco números siguientes:

Mínimo, Q_1 , Mediana, Q_3 , Máximo

Ejemplo: salarios en España

Encuesta de estructura salarial. Año 2006

MEDIAS Y PERCENTILES. Resultados Nacionales

Ganancia media anual por trabajador por sexo, estudios y media y percentiles.

Unidades: euros

	Media	Percentil 10	Percentil 25	Percentil 50	Percentil 75	Percentil 90
Total						
Todos los estudios	19.680,88	8.201,02	11.903,59	15.740,23	23.285,82	34.889,95
Varones						
Todos los estudios	22.051,08	10.608,18	13.483,55	17.204,27	25.671,40	38.620,68
Mujeres						
Todos los estudios	16.245,17	6.258,13	9.682,85	13.506,00	19.722,20	29.323,84

La **(cuasi)desviación típica** es la raíz cuadrada de S^2 :

$$S = \sqrt{\frac{(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n - 1}}$$

S se usa más que S^2 porque mide la dispersión en la misma escala que los datos originales.

Para comparar la dispersión de variables de magnitudes muy distintas a veces se usa el **coeficiente de variación**:

$$CV = \frac{S}{|\bar{x}|}.$$

El CV no depende de las unidades en las que midamos una variable.

Una fórmula alternativa para calcular S^2 :

$$S^2 = \frac{n}{n-1} \left(\frac{x_1^2 + \dots + x_n^2}{n} - \bar{x}^2 \right)$$

Medidas de dispersión: la varianza y la desviación típica

Son las medidas de dispersión más utilizadas.

La **varianza** es el promedio de las desviaciones al cuadrado de los datos a su media.

Datos	x_1, \dots, x_n
Desviaciones	$x_1 - \bar{x}, \dots, x_n - \bar{x}$
Desviaciones al cuadrado	$(x_1 - \bar{x})^2, \dots, (x_n - \bar{x})^2$

La varianza es

$$\frac{(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n}$$

Se suele usar más la **(cuasi)varianza**:

$$S^2 = \frac{(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n - 1}$$

Cuestiones

Da un ejemplo de un conjunto de datos tal que $S^2 = 0$.

Dado un conjunto de observaciones medidas en kg, supongamos que cambiamos las unidades y las pasamos a gramos (es decir, multiplicamos por mil). Determina si son verdaderas o falsas las siguientes afirmaciones:

- ▶ Tanto la media como la mediana de los nuevos datos se multiplican también por mil.
- ▶ La varianza se multiplica también por mil.

¿Cómo cambiaría la desviación típica?

Ahora sumamos 100 a todos los datos. Determina si son verdaderas o falsas las siguientes afirmaciones:

- ▶ Los cuartiles no cambian.
- ▶ El rango intercuartílico no cambia.
- ▶ La desviación típica no cambia.

Descripción numérica

Estadísticos

		LONG	PESO	CONC
N	Válidos	171	171	171
	Perdidos	0	0	0
	Media	39,9708	1147,9123	1,1918
	Error típ. de la media	,65132	66,95359	,05825
	Mediana	39,0000	873,0000	,9300
	Desv. típ.	8,51715	875,53176	,76166
	Varianza	72,542	766555,869	,580
	Rango	39,80	4308,00	3,49
	Mínimo	25,20	203,00	,11
Percentiles	Máximo	65,00	4511,00	3,60
	25	33,3000	491,0000	,5900
	50	39,0000	873,0000	,9300
	75	46,2000	1455,0000	1,6000

Cuestiones

- Calcula el coeficiente de variación de las tres variables. ¿Qué se deduce sobre la dispersión de los valores que toman?
- Comparando los valores de la media y la mediana, ¿cuál de las tres distribuciones parece ser más simétrica?
- Verdadero o falso: Al menos para 100 peces, la concentración de mercurio es superior a 0.93 ppm.
- Verdadero o falso: La longitud de aproximadamente 42 peces es mayor que 25.20 cm y menor que 33.3 cm.
- ¿Cuál es el rango intercuartílico de la variable que mide el peso de los peces?

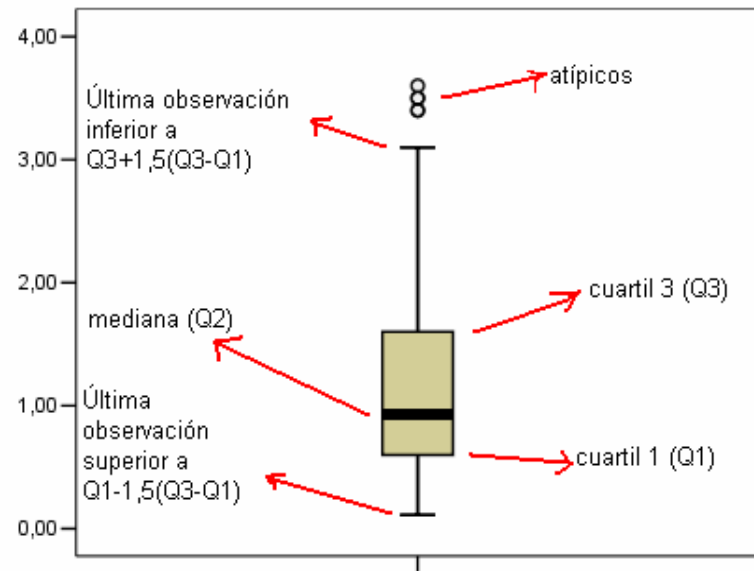
Con SPSS

The screenshot shows the SPSS Statistics Editor window with the file 'mercurio.sav [Conjunto_de_datos1]'. The 'Análisis' menu is open, and the 'Estadísticos descriptivos' submenu is displayed. The 'Frecuencias...' option is highlighted. The data table shows columns for 'RIO', 'ESTACION', 'LONG', 'PESO', and 'CONC'.

Con SPSS

The screenshot shows the SPSS Statistics Editor window with the file 'mercurio.sav [Conjunto_de_datos1]'. The 'Frecuencias' dialog box is open, showing the variables 'LONG', 'PESO', and 'CONC' selected for analysis. The 'Mostrar tablas de frecuencias' checkbox is checked. The data table shows columns for 'RIO', 'ESTACION', 'LONG', 'PESO', and 'CONC'.

Diagrama de cajas



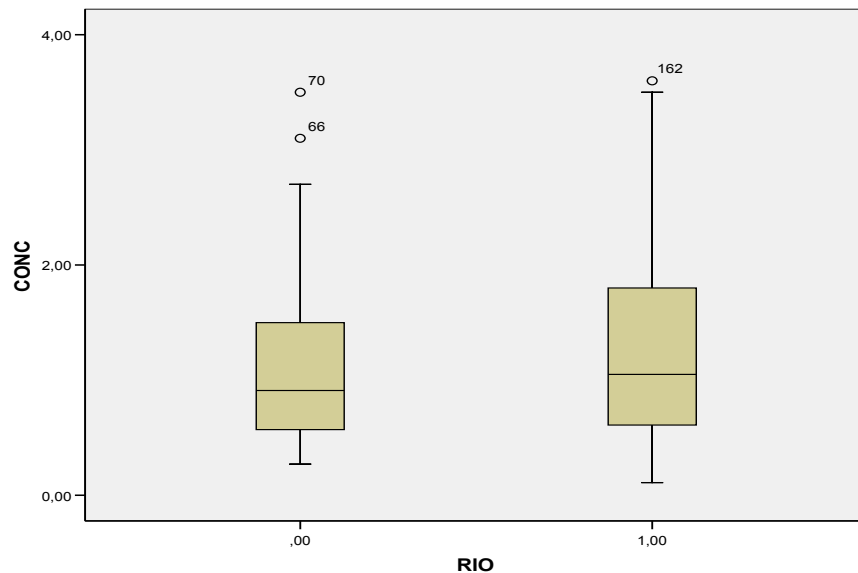
¿Para qué sirven?

Los diagramas de cajas son especialmente útiles para comparar varios conjuntos de datos.

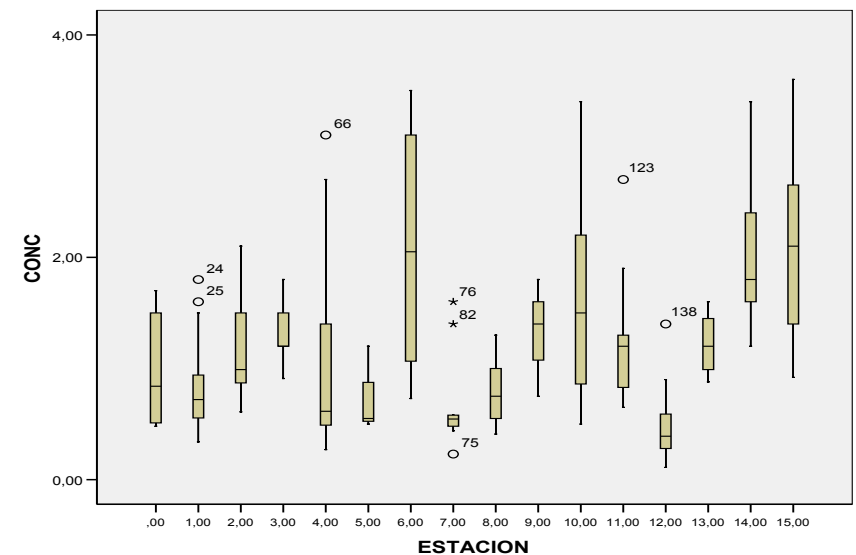
Además, proporcionan información sobre:

- ▶ La posición (mediana) y la dispersión (rango intercuartílico) de los datos.
- ▶ La simetría de la distribución (comparamos el tamaño de las cajas).
- ▶ La existencia de datos que se desvían del patrón general (datos atípicos).

Concentración de mercurio y río



Concentración de mercurio y estación



Relaciona cada histograma con su diagrama de cajas

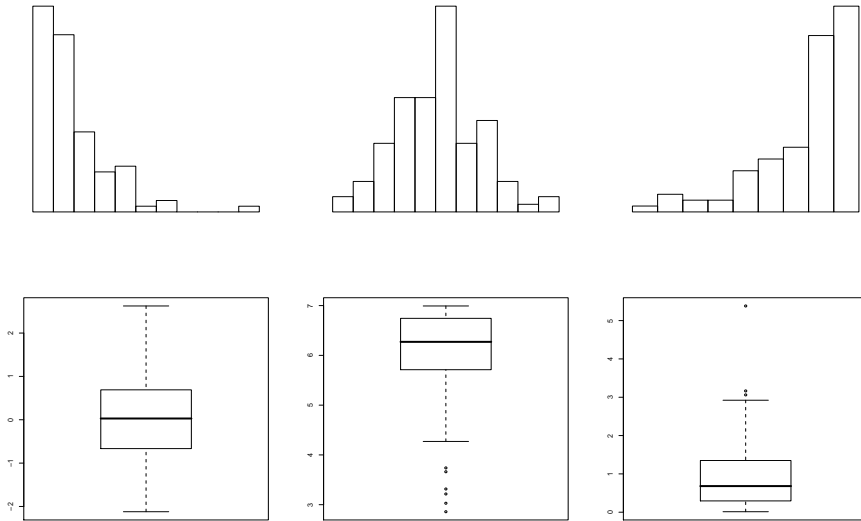
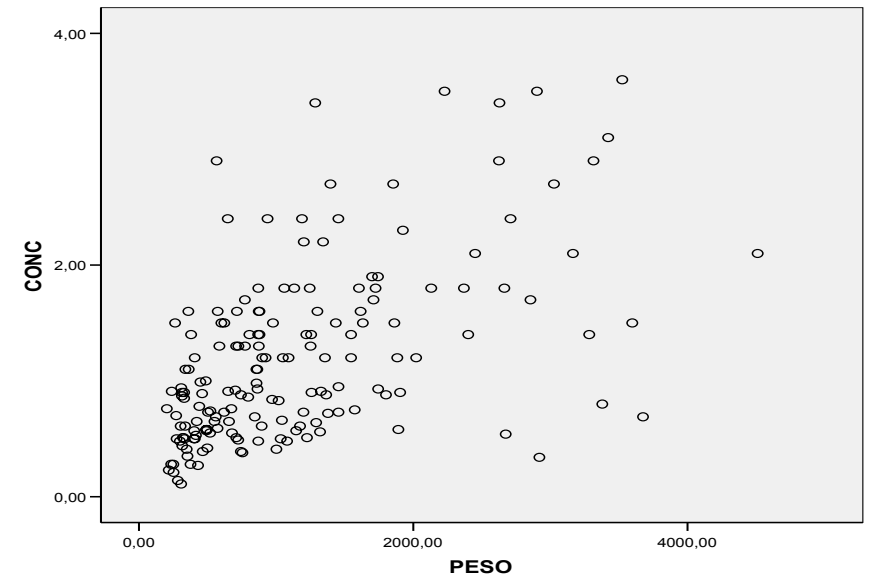


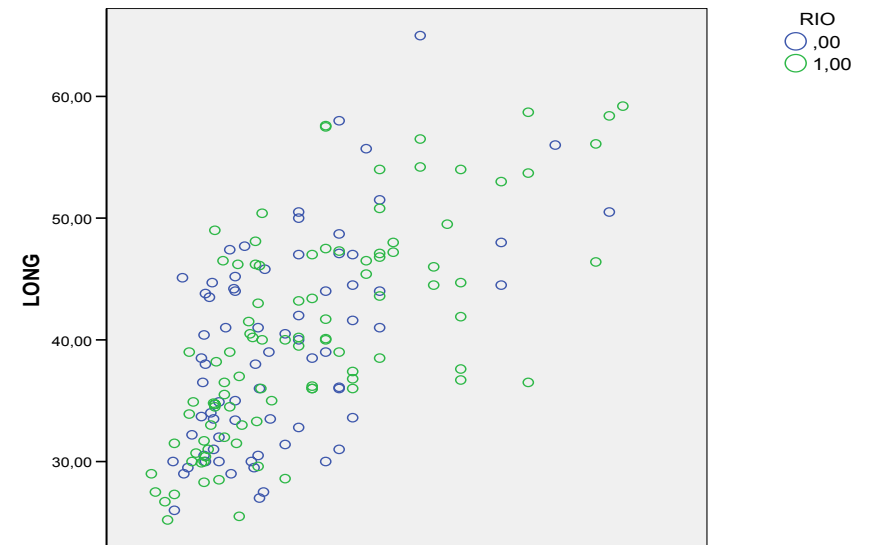
Diagrama de dispersión: Concentración frente a peso



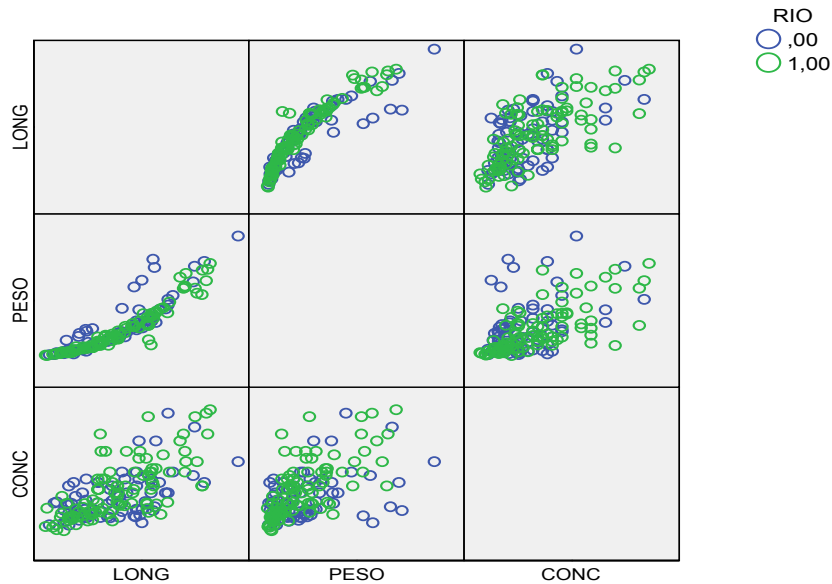
Interpretación de un diagrama de dispersión

- ▶ Es importante fijarse en las unidades de cada eje
- ▶ ¿Se observa alguna asociación entre las variables?
- ▶ ¿Cómo es de estrecha la asociación entre las variables?
- ▶ ¿Cuál es la “dirección” de la asociación entre las variables?
- ▶ ¿Hay algún punto o colección de puntos que no siga el patrón general del resto?
- ▶ Si hay una tercera variable cualitativa, resulta conveniente utilizar símbolos o colores diferentes para cada valor de esta tercera variable.

Concentración frente a longitud (color según río)



Matriz de diagramas de dispersión



Covarianza

Se dispone de un conjunto de n pares de observaciones

$$(x_1, y_1), \dots, (x_n, y_n).$$

El objetivo es definir una medida numérica para cuantificar el grado de relación lineal que hay entre x e y :

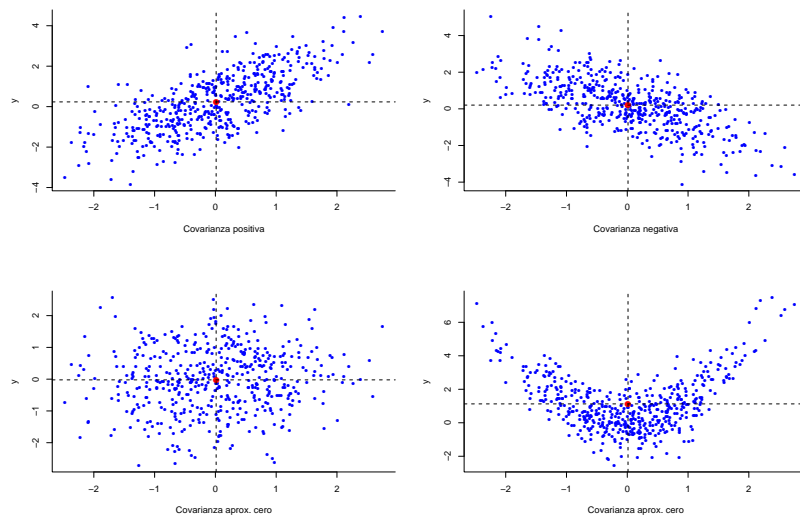
Para ello se usa la covarianza entre x e y :

$$S_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Observaciones:

- ▶ Para entender por qué esta definición es útil miramos el gráfico de la transparencia siguiente.
- ▶ $S_{xy} = S_{yx}$.
- ▶ S_{xx} es la varianza de x .
- ▶ S_{xy} depende de las unidades en que se midan x e y .

Interpretación de la covarianza



Coeficiente de correlación

Resulta conveniente disponer de una medida de relación lineal que no dependa de las unidades. Para ello, se normaliza S_{xy} dividiendo por el producto de desviaciones típicas, lo que lleva al **coeficiente de correlación**:

$$r_{xy} = \frac{S_{xy}}{S_x S_y}.$$

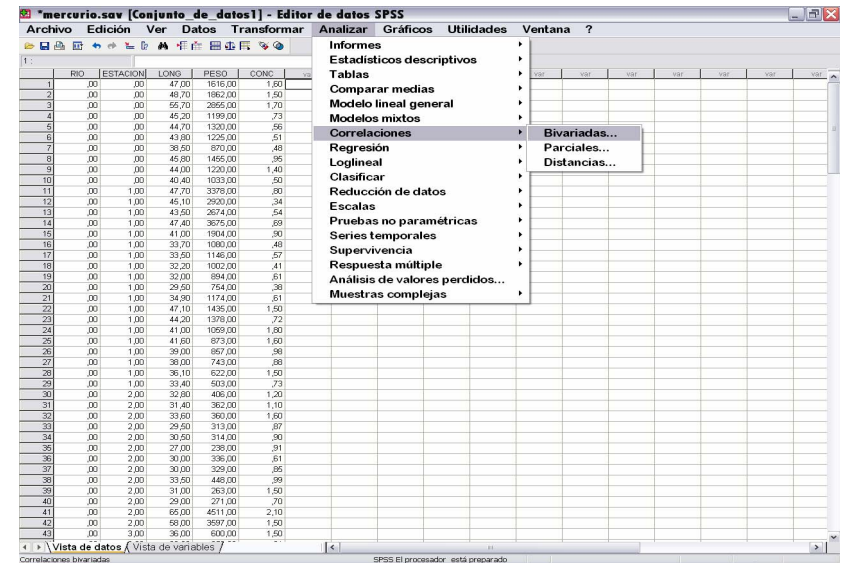
Propiedades:

- ▶ No depende de las unidades
- ▶ Siempre toma valores entre -1 y 1.
- ▶ Su signo se interpreta igual que el de la covarianza
- ▶ Sólo vale 1 ó -1 cuando los puntos están perfectamente alineados.

Covarianzas y correlaciones de los datos

Correlaciones		LONG	PESO	CONC
LONG	Correlación de Pearson	1	,900	,650
	Sig. (bilateral)		,000	,000
	Suma de cuadrados y productos cruzados	12332,114	1141004	716,835
	Covarianza	72,542	6711,790	4,217
	N	171	171	171
PESO	Correlación de Pearson	,900	1	,554
	Sig. (bilateral)	,000		,000
	Suma de cuadrados y productos cruzados	1141004	1E+008	62786,546
	Covarianza	6711,790	766555,9	369,333
	N	171	171	171
CONC	Correlación de Pearson	,650	,554	1
	Sig. (bilateral)	,000	,000	
	Suma de cuadrados y productos cruzados	716,835	62786,546	98,622
	Covarianza	4,217	369,333	,580
	N	171	171	171

Covarianzas y correlaciones con SPSS



Estandarización o tipificación

Consiste en restarle a cada observación la media de todos los datos y dividir por la desviación típica:

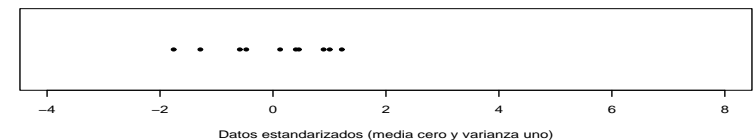
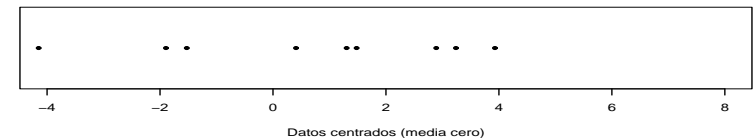
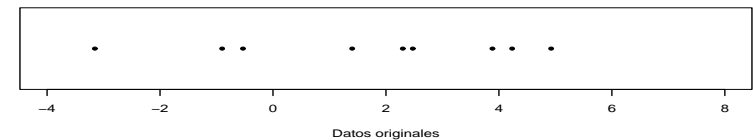
$$z_i = \frac{x_i - \bar{x}}{S}$$

Representa la distancia de x_i a la media expresada en desviaciones típicas (el signo indica si el dato es mayor o menor que la media).

¿Cuánto vale la media de los datos estandarizados?

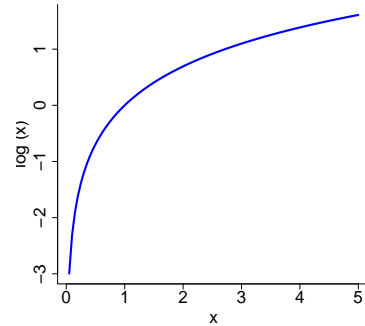
¿Y su desviación típica?

Efecto de estandarizar un conjunto de datos

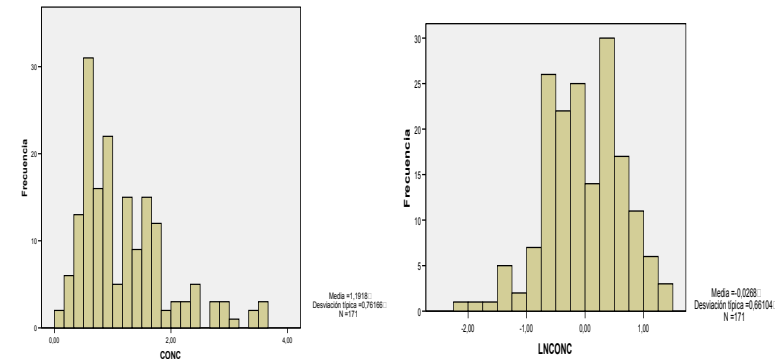


Tomar logaritmos

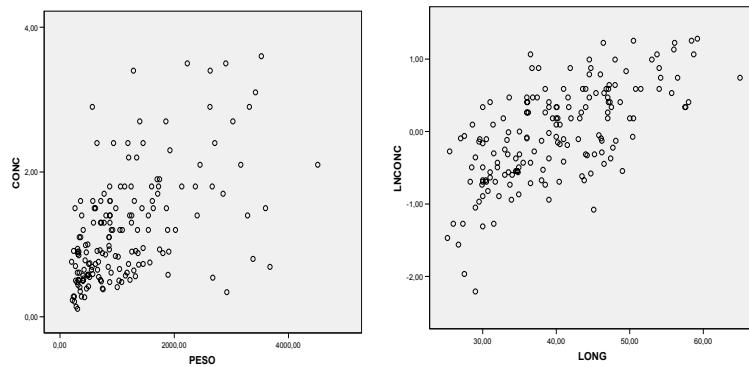
Si las observaciones x_i son positivas, a veces es conveniente trabajar con sus logaritmos $\log x_i$ en lugar de con las variables originales.



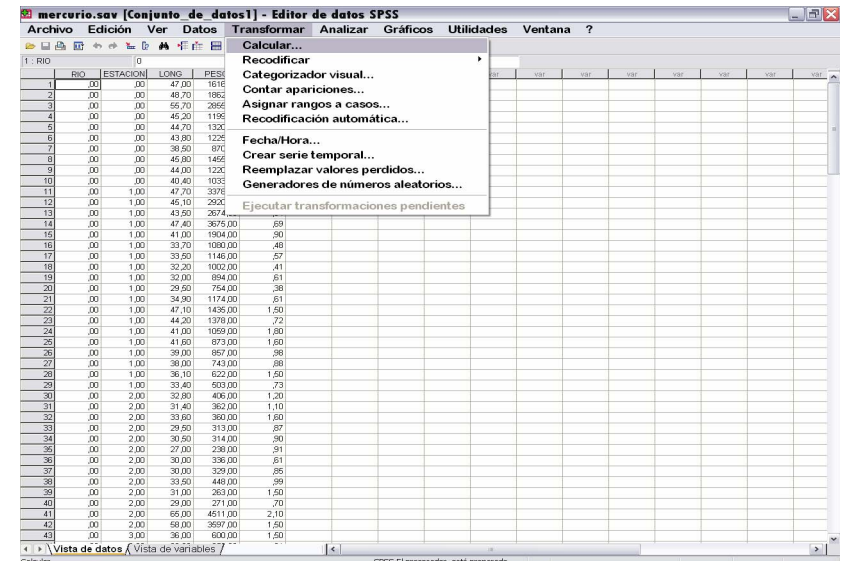
Tomar logaritmos para hacer la distribución más simétrica



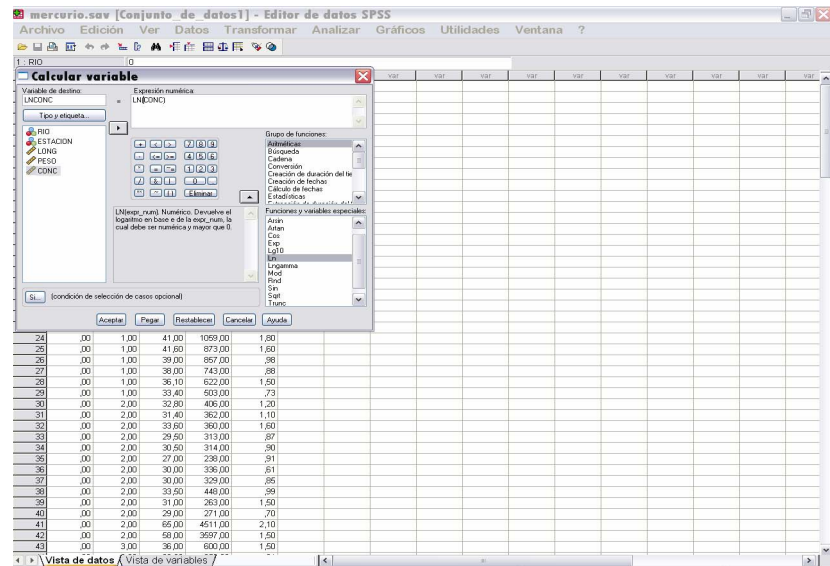
Tomar logaritmos para hacer que la asociación sea lineal



Transformaciones con SPSS



Transformaciones con SPSS



Calorías y contenido en sodio en salchichas

- Se ha considerado la cantidad de calorías y de sodio en salchichas de varias marcas de cada uno de los tipos siguientes:

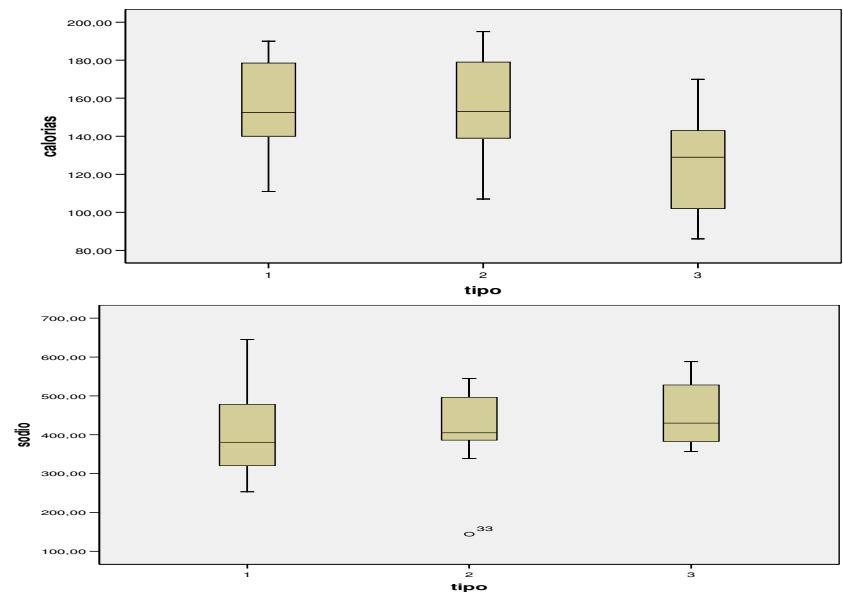
- Carne de ternera
- Mezcla (hasta 15% de carne de pavo)
- Carne de pavo

Nombre variable	Descripción
tipo	Tipo de carne (1=ternera, 2=mezcla, 3=pavo)
calorias	Cantidad de calorías
sodio	Cantidad de sodio

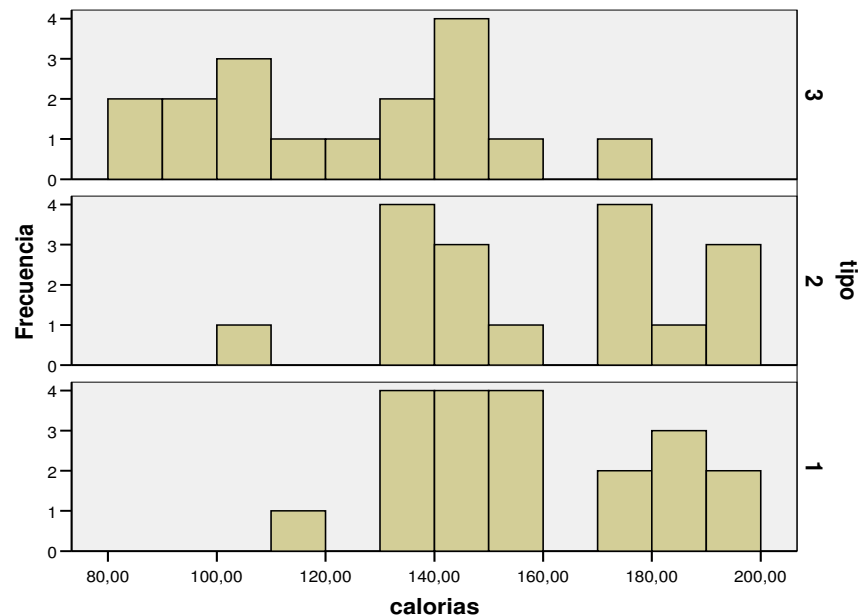
Medidas descriptivas numéricas

		calorias	sodio
N	Válidos	54	54
	Perdidos	0	0
Media		146,6111	424,8333
Error típ. de la media		3,95691	13,04440
Mediana		146,0000	405,0000
Desv. típ.		29,07727	95,85637
Varianza		845,487	9188,443
Mínimo		86,00	144,00
Máximo		195,00	645,00
Percentiles	25	132,0000	359,7500
	50	146,0000	405,0000
	75	173,5000	506,2500

Diagramas de cajas



Histogramas: cantidad de calorías



Histogramas: cantidad de sodio

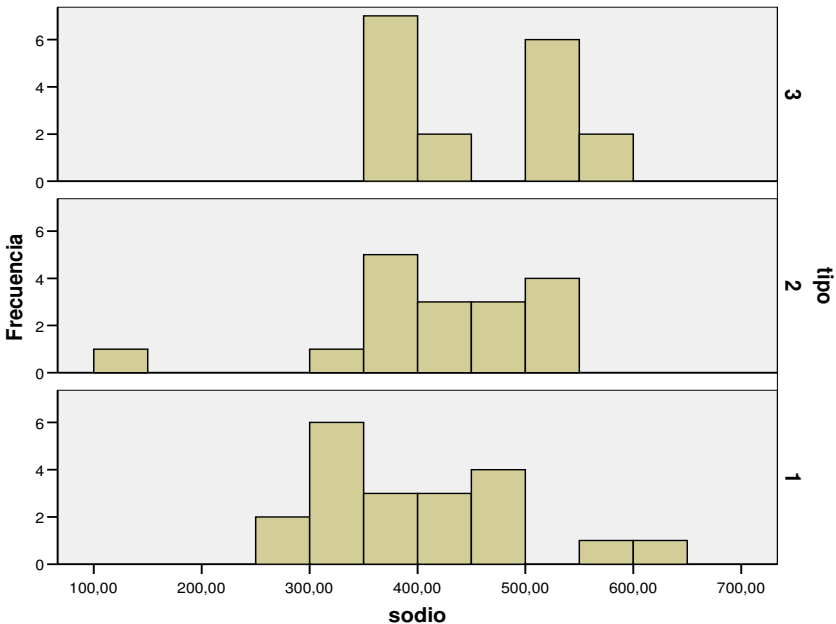
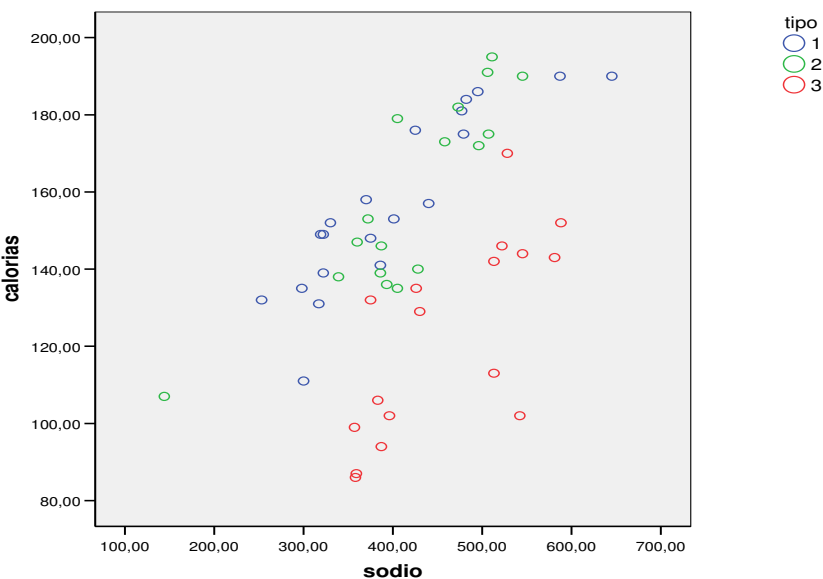


Diagrama de dispersión



Covarianzas y correlaciones

Correlaciones

		calorias	sodio
calorias	Correlación de Pearson	1	,516
	Sig. (bilateral)		,000
	Suma de cuadrados y productos cruzados	44810,833	76233,500
	Covarianza	845,487	1438,368
	N	54	54
sodio	Correlación de Pearson	,516	1
	Sig. (bilateral)	,000	
	Suma de cuadrados y productos cruzados	76233,500	486987,50
	Covarianza	1438,368	9188,443
	N	54	54

Cuestiones

- ▶ (V ó F) Aproximadamente 27 marcas de salchichas tienen entre 132 y 173 calorías.
- ▶ ¿Cuál es el rango intercuartílico de la cantidad de sodio?
- ▶ Calcula el coeficiente de variación de ambas variables.
- ▶ (V ó F) Aproximadamente 13 marcas de salchichas tienen un contenido de sodio entre 506.25 y 645.
- ▶ (V ó F) Con la información disponible en la tabla de medidas descriptivas numéricas es posible calcular la correlación entre ambas variables.
- ▶ (V ó F) Al menos el 75% de las marcas de salchichas de mezcla tienen más sodio que la mediana de las marcas de ternera.
- ▶ Identifica en el diagrama de dispersión el dato atípico que se observa en los diagramas de cajas.