

NOTAS DE CLASE :

**Introducción a la Estadística
con Apoyo Computacional**

Felipe Osorio y Ronny Vallejos

DEPARTAMENTO DE MATEMÁTICA, UNIVERSIDAD TÉCNICA FEDERICO SANTA MARÍA

Índice general

Capítulo 1. Estadística Descriptiva	1
1.1. Preliminares: Sumas y Productos	1
1.2. Estadísticas de Resumen	5
1.3. Covarianza y correlación	27
1.4. Estadísticas descriptivas multivariadas	30
1.5. Regresión lineal simple	34
1.6. Resúmenes gráficos	39
Capítulo 2. Nociones de Probabilidad	45
2.1. Preliminares	45
2.2. Conceptos básicos	46
2.3. Espacios muestrales finitos	49
2.4. Técnicas de conteo	50
2.5. Probabilidad condicional	53
2.6. Independencia estadística	55
Bibliografía	57

Estadística Descriptiva

1.1. Preliminares: Sumas y Productos

En esta sección se introduce notación que tiene por objetivo escribir de forma compacta sumas y productos de secuencias de números a_1, a_2, \dots , donde $a_i \in \mathbb{R}$, para todo i .

DEFINICIÓN 1.1. Considere una secuencia de números a_1, a_2, \dots . Se define la *sumatoria* de esta secuencia, como:

$$\sum_{i=1}^n a_i = a_1 + a_2 + \dots + a_n, \quad (1.1)$$

donde i denota el índice de la sumatoria, mientras que a_i representa un elemento genérico. En este caso, n indica la cantidad de elementos que se están sumando.

Es posible apreciar que la suma en (1.1) puede ser escrita de manera análoga como

$$\sum_{1 \leq i \leq n} a_i = a_1 + a_2 + \dots + a_n. \quad (1.2)$$

Debemos notar que, si $n = 0$ el valor de la sumatoria se define como cero.

OBSERVACIÓN 1.2. A partir de la Ecuación (1.2) podemos introducir una notación mucho más general. En efecto, sea R un conjunto de índices. Así, basta considerar el conjunto $R = \{1, 2, \dots, n\}$, para re-escribir la suma en (1.2) como:

$$\sum_{i \in R} a_i = a_1 + \dots + a_n. \quad (1.3)$$

OBSERVACIÓN 1.3. Aunque frecuentemente la notación dada en la Ecuación (1.3) es utilizada para sumas finitas, esta puede ser adaptada con facilidad para sumas infinitas. Por ejemplo,

$$\sum_{i=1}^{\infty} a_i = \sum_{i \geq 1} a_i = a_1 + a_2 + \dots.$$

Más formalmente, debemos escribir

$$\sum_{i=1}^{\infty} a_i = \lim_{n \rightarrow \infty} \sum_{i=1}^n a_i.$$

RESULTADO 1.4 (Regla constante). Sea a un número real. De este modo,

$$\sum_{i=1}^n a = a + a + \dots + a = na.$$

En general, para $r < n$ tenemos

$$\sum_{i=r}^n a = (n - r + 1)a, \quad a \in \mathbb{R}.$$

RESULTADO 1.5. Considere la secuencia x_1, \dots, x_n y sea a una constante. Entonces,

$$\sum_{i=1}^n a x_i = a x_1 + \dots + a x_n = a(x_1 + \dots + x_n) = a \sum_{i=1}^n x_i.$$

En general, sean x_1, \dots, x_n y y_1, \dots, y_n dos secuencias de números y $a, b \in \mathbb{R}$. Entonces,

$$\sum_{i=1}^n (a x_i + b y_i) = a \sum_{i=1}^n x_i + b \sum_{i=1}^n y_i.$$

Note también que las sumatorias pueden ser *descompuestas* en varias sumas. En efecto, para una secuencia de números a_1, \dots, a_n . Tenemos que

$$\sum_{i=1}^n a_i = \sum_{i=1}^k a_i + \sum_{i=k+1}^n a_i, \quad k < n.$$

En general, sea $R = R_1 \cup R_2$, tal que $R_1 \cap R_2 = \emptyset$. Entonces,

$$\sum_{i \in R} a_i = \sum_{i \in R_1} a_i + \sum_{i \in R_2} a_i.$$

EJEMPLO 1.6 (Propiedad telescópica). Suponga a_0, a_1, \dots, a_n una secuencia de números reales, y considere

$$\begin{aligned} \sum_{i=1}^n (a_i - a_{i-1}) &= (a_1 - a_0) + (a_2 - a_1) + \dots + (a_{n-1} - a_{n-2}) + (a_n - a_{n-1}) \\ &= -a_0 + (a_1 - a_1) + \dots + (a_{n-1} - a_{n-1}) + a_n \\ &= a_n - a_0. \end{aligned}$$

EJEMPLO 1.7 (Suma de una progresión geométrica). Asuma que $x \neq 1$ y $n \geq 0$. Entonces,

$$\begin{aligned} a + ax + \dots + ax^n &= \sum_{j=0}^n ax^j \\ &= a + \sum_{j=1}^n ax^j = a + (ax + \dots + ax^n) = a + x(a + \dots + ax^{n-1}) \end{aligned} \quad (1.4)$$

$$= a + x \sum_{j=1}^n ax^{j-1} = a + x \sum_{j=0}^{n-1} ax^j \quad (1.5)$$

$$= a + x \sum_{j=0}^n ax^j - ax^{n+1}. \quad (1.6)$$

Tomando la primera y última relaciones, sigue que

$$\sum_{j=0}^n ax^j = a + x \sum_{j=0}^n ax^j - ax^{n+1},$$

es decir,

$$(1-x) \sum_{j=0}^{\infty} ax^j = a - ax^{n+1}.$$

De este modo, obtenemos finalmente

$$\sum_{j=0}^n ax^j = a \left(\frac{1-x^{n+1}}{1-x} \right).$$

En este ejemplo se han utilizado diversos elementos que permiten notar algunas de las propiedades de las sumatorias. En efecto, en Ecuación (1.4) se utilizó una muy particular versión del Resultado 1.5. Mientras que en (1.5) y (1.6) se realizó un cambio en el índice de la suma y reorganizó los términos (es decir el dominio sobre el que opera la suma) para obtener el resultado deseado.

Las siguientes son igualdades que **no** **satisface** la suma:

- Sean a_1, \dots, a_n y b_1, \dots, b_n dos secuencias de números reales. Entonces

$$\sum_{i=1}^n a_i b_i \neq \left(\sum_{i=1}^n a_i \right) \left(\sum_{i=1}^n b_i \right). \quad (1.7)$$

En efecto, basta notar que la cantidad de términos involucrados en cada uno de los lados de la ecuación anterior es diferente.

- Un caso particular del anterior es

$$\sum_{i=1}^n x_i^2 \neq \left(\sum_{i=1}^n x_i \right)^2.$$

- En general, si $f : \mathbb{R} \rightarrow \mathbb{R}$ es una función no lineal. Entonces

$$\sum_{i=1}^n f(x_i) \neq f \left(\sum_{i=1}^n x_i \right).$$

En ocasiones disponemos de secuencias de números indexados mediante dos (o más) índices, es decir $\{a_{ij}\}$. Por ejemplo, consideremos $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{m \times n}$ (es decir $i = 1, \dots, m$, $j = 1, \dots, n$) y suponga que deseamos sumar todos los elementos de la matriz \mathbf{A} . Es decir,

$$\sum_{i=1}^m \sum_{j=1}^n a_{ij} = a_{11} + \dots + a_{1n} + a_{21} + \dots + a_{m1} + \dots + a_{mn}.$$

Notamos fácilmente que podemos intercambiar el orden de las sumas. En efecto,

$$\sum_{i=1}^m \sum_{j=1}^n a_{ij} = \sum_{j=1}^n \sum_{i=1}^m a_{ij}$$

OBSERVACIÓN 1.8. Se debe resaltar que la operación de intercambiar el orden de las sumas **no** siempre es válido para series infinitas.

Contrariamente al resultado de la Ecuación (1.7), es válido considerar

$$\left(\sum_{i=1}^m a_i \right) \left(\sum_{j=1}^n b_j \right) = \sum_{i=1}^m \sum_{j=1}^n a_i b_j, \quad (1.8)$$

asimismo

$$\left(\sum_{i=1}^n x_i\right)^2 = \left(\sum_{i=1}^n x_i\right)\left(\sum_{j=1}^n x_j\right) = \sum_{i=1}^n \sum_{j=1}^n x_i x_j.$$

Para comprender mejor la Ecuación (1.8) considere un caso especial

$$\begin{aligned} \left(\sum_{i=1}^2 a_i\right)\left(\sum_{j=1}^3 b_j\right) &= (a_1 + a_2)(b_1 + b_2 + b_3) \\ &= (a_1 b_1 + a_1 b_2 + a_1 b_3) + (a_2 b_1 + a_2 b_2 + a_2 b_3) \\ &= \sum_{i=1}^2 \left(\sum_{j=1}^3 a_i b_j\right). \end{aligned}$$

Otros ejemplos de sumas útiles (que pueden ser probadas usando inducción) son:

- $\sum_{k=1}^n k = 1 + 2 + \cdots + n = \frac{n(n+1)}{2}.$
- $\sum_{k=1}^n k^2 = 1^2 + 2^2 + \cdots + n^2 = \frac{n(n+1)(2n+1)}{6}.$

Existe una notación análoga para productos. En efecto, considere la siguiente definición

DEFINICIÓN 1.9. Sea a_1, a_2, \dots una secuencia de números. Se define la *productoria* de esta secuencia, como:

$$\prod_{i=1}^n a_i = a_1 a_2 \cdots a_n. \quad (1.9)$$

En general, podemos escribir

$$\prod_{i \in R} a_i,$$

donde R representa un conjunto de índices. Note que si no existe algún entero $i \in R$, el producto se define con el valor uno.

EJEMPLO 1.10 (factorial de un número). Un ejemplo del uso de productorios es:

$$1 \cdot 2 \cdot 3 \cdots (n-1) \cdot n = \prod_{j=1}^n j = n!$$

que se denomina n factorial. Recuerde que $0!$ por definición es 1.

OBSERVACIÓN 1.11. En efecto, el factorial tiene una propiedad interesante. Es muy fácil notar que:

$$n! = (n-1)!n,$$

lo que permite definir el factorial de manera *recursiva*.

1.2. Estadísticas de Resumen

1.2.1. Medidas de Posición. Se introduce una serie de medidas que permiten resumir un gran volumen de información y cuantifican el valor central de un conjunto de datos.

DEFINICIÓN 1.12 (Media muestral o promedio). Sea x_1, \dots, x_n valores muestrales. Se define el *promedio* o *media muestral* como:

$$\bar{x} = \frac{1}{n}(x_1 + \dots + x_n) = \frac{1}{n} \sum_{i=1}^n x_i. \quad (1.10)$$

Suponga que la observación i -ésima, digamos x_i , se repite n_i veces. Entonces tenemos que

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n n_i x_i = \sum_{i=1}^n f_i x_i, \quad (1.11)$$

donde $f_i = n_i/n$ es la frecuencia relativa. Considere “pesos” o ponderaciones $\omega_1, \dots, \omega_n$ asociados a las observaciones x_1, \dots, x_n . En este caso tenemos,

$$\bar{x} = \frac{1}{\sum_{j=1}^n \omega_j} \sum_{i=1}^n \omega_i x_i, \quad (1.12)$$

si consideramos la proporción

$$p_i = \frac{\omega_i}{\sum_{j=1}^n \omega_j}, \quad i = 1, \dots, n,$$

entonces podemos re-escribir (1.12) como

$$\bar{x} = \sum_{i=1}^n p_i x_i.$$

Note que el promedio en Ecuación (1.11) es un caso particular donde $\sum_i f_i = 1$.

EJEMPLO 1.13. Considere el conjunto de datos $\{1, 2, 2, 2, 3, 3, 8\}$. Tenemos $n = 7$, y

$$\sum_{i=1}^7 x_i = 1 + 3 \cdot 2 + 2 \cdot 3 + 8 = 21,$$

así $\bar{x} = 21/7 = 3$. Note también que el gráfico de *tallo y hoja*, adopta la forma:

1	*			
2	*	*	*	*
3	*	*		
4				
5				
6				
7				
8	*			

EJEMPLO 1.14 (Datos de accidentes). Suponga el siguiente conjunto de datos:

Número de accidentes (x_i)	Frecuencia (n_i)	$n_i x_i$
0	55	0
1	14	14
2	5	10
3	2	6
4	0	0
Total	76	30

De este modo, $\bar{x} = 30/76 = 0.395$ es el número promedio de accidentes.

Considere que el conjunto de observaciones x_1, \dots, x_n es ordenado de menor a mayor como:

$$x_{(1)}, x_{(2)}, \dots, x_{(n)},$$

tal que $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$, donde $x_{(k)}$ se denomina la k -ésima estadística de orden.

DEFINICIÓN 1.15 (Mediana). Sean $x_{(1)}, \dots, x_{(n)}$ observaciones ordenadas. Si n es impar, entonces la *mediana* se define como la observación central, es decir

$$\mathbf{me} = x_{((n+1)/2)},$$

en el caso que n sea par, entonces

$$\mathbf{me} = \frac{1}{2} \left(x_{(n/2)} + x_{(n/2+1)} \right).$$

OBSERVACIÓN 1.16. En ocasiones escribiremos $\mathbf{me}(\mathbf{x})$ para indicar cual es el conjunto de datos sobre el que se calcula la mediana.

DEFINICIÓN 1.17 (Moda). La *moda* o *valor modal* es el valor observado con la más alta ocurrencia.

OBSERVACIÓN 1.18. Respecto de las medidas de tendencia central introducidas anteriormente podemos apreciar que:

- El promedio puede verse fuertemente afectado por *datos atípicos*.
- La mediana “divide” el conjunto de datos en dos, es decir, el 50 % de los datos están por debajo de la mediana, mientras que el 50 % se encuentran por sobre este valor.
- En general, la moda no es única y puede no existir.
- Es interesante notar la diferencia entre la complejidad de cálculo del promedio versus el de la mediana.

OBSERVACIÓN 1.19 (Otras medidas de tendencia central). Sea $f(x)$ cualquier función de números reales. Entonces podemos definir

$$\bar{f} = \frac{1}{n} \sum_{i=1}^n f(x_i) = \frac{1}{n} \left(f(x_1) + \dots + f(x_n) \right).$$

Los siguientes casos particulares son de interés:

- (a) *Média cuadrática*. Considere $f(x) = x^2$. Entonces se define la media cuadrática, Q como:

$$Q = \frac{1}{n} \sum_{i=1}^n x_i^2.$$

- (b) *Média armónica*. Sea $f(x) = 1/x$. Entonces, decimos que H es la media armónica si

$$\frac{1}{H} = \frac{1}{n} \left(\frac{1}{x_1} + \frac{1}{x_2} + \cdots + \frac{1}{x_n} \right) = \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}.$$

Es decir, H es el inverso de la media aritmética de los inversos de los valores observados. De donde sigue que,

$$H = \frac{n}{\sum_{i=1}^n 1/x_i}.$$

- (c) *Média geométrica*. Considere $f(x) = \log x$. Entonces la media geométrica G es definida por la fórmula

$$\log G = \frac{1}{n} \left(\log x_1 + \cdots + \log x_n \right) = \frac{1}{n} \sum_{i=1}^n \log x_i. \quad (1.13)$$

Es decir,

$$G = \left(\prod_{i=1}^n x_i \right)^{1/n}.$$

OBSERVACIÓN 1.20. Un procedimiento bastante usado para el cálculo de la média geométrica es obtener la média aritmética dada en la Ecuación (1.13), y luego considerar

$$G = \exp \left(\frac{1}{n} \sum_{i=1}^n \log x_i \right).$$

Sin embargo, este método aunque correcto, puede inducir algunos errores en la precisión de los resultados. Una alternativa que intenta corregir esta situación, es basada en el trabajo de [Graillat \(2009\)](#), quien propuso un esquema compensado para la evaluación de productos utilizando un mecanismo de multiplicación y suma fundidas (FMA). La biblioteca `fastmatrix` ([Osorio y Ogeda, 2022](#)), contiene una implementación de la média geométrica usando este enfoque. El siguiente ejemplo, con comandos en R ilustra el uso de la función `geomean`.

```
# introduciendo datos en la consola de R
> x <- c(2.2, 0.3, 0.5, 0.4, 0.2, 1.9)

# cargando biblioteca 'fastmatrix'
> library(fastmatrix)
> geomean(x)
[1] 0.6072855
> exp(mean(log(x))) # equivalente a 'geomean'
[1] 0.6072855
```

1.2.2. Medidas de Dispersión. Considere los conjuntos de datos:

$$D_1 = \{10, 20, 30\}, \quad D_2 = \{5, 5, 20, 35, 35\}, \quad D_3 = \{20, 20, 20\}.$$

Tenemos los gráficos de tallo-y-hoja:

Datos D_1 :	Datos D_2 :	Datos D_3 :
5	5 * *	5
10 *	10	10
15	15	15
20 *	20 *	20 * * *
25	25	25
30 *	30	30
35	35 * *	35

Sea \bar{x}_j y me_j el promedio y la mediana asociada al conjunto de datos D_j ($j = 1, 2, 3$). Entonces,

$$\begin{aligned} \bar{x}_1 &= \frac{1}{3}(10 + 20 + 30) = \frac{60}{3} = 20, \\ \bar{x}_2 &= \frac{1}{5}(2 \cdot 5 + 20 + 2 \cdot 35) = \frac{100}{5} = 20, \\ \bar{x}_3 &= \frac{3 \cdot 20}{3} = 20. \end{aligned}$$

Además, $\text{me}_j = 20$ para $j = 1, 2, 3$. Es decir, tenemos tres configuraciones de datos con valores centrales idénticos.

Esto motiva la introducción de medidas que permitan caracterizar la variabilidad o dispersión de un conjunto de observaciones. Algunas medidas simples corresponden a los *cuartiles*. En efecto, sea Q_1 y Q_3 las medianas de la mitad inferior y superior de los datos, respectivamente. Entonces, esto lleva a definir la siguiente medida de dispersión:

$$IQR = Q_3 - Q_1,$$

el que es conocido como *rango intercuartílico*. Es interesante notar que algunos software estadísticos (por ejemplo, R/S-Plus, Stata, entre otros) reportan:

$$\min\{x_i\}_{i=1}^n, Q_1, \text{me}, Q_3, \max\{x_i\}_{i=1}^n.$$

Así, también podemos considerar el *rango* de la muestra como

$$R = \max\{x_i\}_{i=1}^n - \min\{x_i\}_{i=1}^n = x_{(n)} - x_{(1)}.$$

Por otro lado, se ha sugerido utilizar subdivisiones más finas que los cuartiles. Por ejemplo, considere subdividir los datos ordenados $x_{(1)}, x_{(2)}, \dots, x_{(n)}$ en secciones de 100%, llamados *percentiles*. Precisamente, si disponemos de los datos ordenados $x_{(1)}, x_{(2)}, \dots, x_{(n)}$, entonces el percentil de orden j ($1 \leq j \leq 100$) está dado por:

$$P_j = x_{(j(n+1)/100)}. \quad (1.14)$$

Debemos notar además que, el primer cuartil Q_1 corresponde al percentil 25^o, mientras que la mediana (o 2^o cuartil, Q_2) representa el percentil 50^o y Q_3 corresponde al percentil 75^o.

EJERCICIO 1.21. Considere el conjunto de datos $\mathbf{x} = (4, 7, 18, 1, 7, 13, 2)^\top$ y suponga que deseamos calcular P_{75} y el rango intercuartílico IQR . Primeramente es necesario ordenar el conjunto de datos:

$$(x_{(1)}, x_{(2)}, x_{(3)}, x_{(4)}, x_{(5)}, x_{(6)}, x_{(7)})^\top = (1, 2, 4, 7, 7, 13, 18)^\top.$$

Disponemos de $n = 7$ datos, luego para obtener el 1er y 3er cuartiles podemos usar la fórmula del percentil en (1.14). En efecto,

$$\begin{aligned} Q_1 = P_{25} &= x_{(25 \cdot (7+1)/100)} = x_{(1.8/4)} = x_{(2)} = 2, \\ Q_3 = P_{75} &= x_{(75 \cdot (7+1)/100)} = x_{(3.8/4)} = x_{(6)} = 13. \end{aligned}$$

De este modo, $IQR = Q_3 - Q_1 = 13 - 2 = 11$.

DEFINICIÓN 1.22 (Varianza muestral). Considere x_1, x_2, \dots, x_n valores observados, se define su *varianza muestral* como:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (1.15)$$

Note que s^2 corresponde a un “promedio” de los desvios al cuadrado con relación a la media. También se suele anotar s_x^2 o bien $\text{var}(x)$.

Aunque en ocasiones es recomendable dividir la suma de cuadrados en (1.15) por n . En el capítulo sobre inferencia veremos la razón de utilizar $n - 1$.

OBSERVACIÓN 1.23. $s = \sqrt{s^2}$ se denomina desviación estándar.

Basados en Ecuación (1.15) podemos definir otras medidas de dispersión:

- *Desviación absoluta en torno de la media*

$$\frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|. \quad (1.16)$$

- *Desviación absoluta en torno de la mediana*

$$\frac{1}{n} \sum_{i=1}^n |x_i - \text{me}|. \quad (1.17)$$

OBSERVACIÓN 1.24. En general podemos considerar, por ejemplo:

$$g(T) = \frac{1}{n} \sum_{i=1}^n h(x_i - T(\mathbf{x})),$$

donde $T(\mathbf{x})$ es alguna estadística de la muestra $\mathbf{x} = (x_1, \dots, x_n)^\top$. Note que si $T(\mathbf{x}) = \bar{x}$ y $h(z) = z^2$, obtenemos la varianza. Mientras que para $T(\mathbf{x}) = \text{me}$ y $h(z) = |z|$ obtenemos (1.17).

De este modo, podemos definir el r -ésimo momento centrado en torno de a , como

$$M_r(a) = \frac{1}{n} \sum_{i=1}^n (x_i - a)^r. \quad (1.18)$$

PROPIEDAD 1.25. A continuación se describen una serie de propiedades del promedio y la varianza de un conjunto de datos observados.

- (a) $\sum_{i=1}^n (x_i - \bar{x}) = 0$. En efecto,

$$\sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} = \sum_{i=1}^n x_i - n\bar{x} = n\bar{x} - n\bar{x} = 0.$$

(b) (Fórmula de Kőning)

$$\begin{aligned}\sum_{i=1}^n (x_i - \bar{x})^2 &= \sum_{i=1}^n (x_i^2 - 2x_i\bar{x} + \bar{x}^2) = \sum_{i=1}^n x_i^2 - 2\bar{x} \sum_{i=1}^n x_i + n\bar{x}^2 \\ &= \sum_{i=1}^n x_i^2 - n\bar{x}^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2\end{aligned}\quad (1.19)$$

(c) \bar{x} es el valor que minimiza la función $S(a) = \sum_{i=1}^n (x_i - a)^2$. En efecto, note que

$$\frac{d}{da} S(a) = \sum_{i=1}^n \frac{d}{da} (x_i - a)^2 = -2 \sum_{i=1}^n (x_i - a),$$

resolviendo la condición de primer orden, tenemos

$$\sum_{i=1}^n (x_i - \hat{a}) = 0,$$

desde donde sigue que $\hat{a} = \bar{x}$. Además

$$\frac{d^2}{da^2} S(a) = -2 \sum_{i=1}^n \frac{d}{da} (x_i - a) = 2n,$$

y como la segunda derivada es positiva (para cualquier valor de n), obtenemos que \bar{x} es mínimo global.

Una manera alternativa para probar este resultado puede ser obtenida mediante notar que

$$\begin{aligned}S(a) &= \sum_{i=1}^n (x_i - a)^2 = \sum_{i=1}^n \{(x_i - \bar{x}) + (\bar{x} - a)\}^2 \\ &= \sum_{i=1}^n (x_i - \bar{x})^2 + 2(\bar{x} - a) \sum_{i=1}^n (x_i - \bar{x}) + \sum_{i=1}^n (\bar{x} - a)^2.\end{aligned}$$

Por la propiedad en (a) y notando que el término $(\bar{x} - a)^2$ es constante para la suma. Obtenemos

$$S(a) = \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - a)^2,$$

y de este modo, resulta evidente que $S(a)$ alcanza su mínimo para $\hat{a} = \bar{x}$.

(d) Sea x_1, x_2, \dots, x_n y considere la transformación

$$y_i = ax_i + b.$$

Entonces

$$\bar{y} = a\bar{x} + b, \quad s_y^2 = a^2 s_x^2.$$

Es fácil notar que,

$$\begin{aligned}\bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (ax_i + b) = \frac{1}{n} \left(a \sum_{i=1}^n x_i + b \right) \\ &= a \left(\frac{1}{n} \sum_{i=1}^n x_i \right) + b = a\bar{x} + b.\end{aligned}$$

Mientras que

$$s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2,$$

como $y_i - \bar{y} = ax_i + b - (a\bar{x} + b) = a(x_i - \bar{x})$, sigue que

$$\begin{aligned} s_y^2 &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n-1} \sum_{i=1}^n \{a(x_i - \bar{x})\}^2 \\ &= a^2 \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = a^2 s_x^2. \end{aligned}$$

En particular, si (que corresponde a una *estandarización* del conjunto de datos x_1, \dots, x_n)

$$z_i = \frac{x_i - \bar{x}}{s}, \quad i = 1, \dots, n.$$

Entonces $\bar{z} = 0$ y $s_z^2 = 1$. En efecto, en la Propiedad 1.25(d) basta hacer $a = 1/s$ y $b = \bar{x}/s$.

- (e) Sea x_1, \dots, x_n un conjunto de n observaciones. Considere aplicar la transformación

$$y_i = g(x_i), \quad i = 1, \dots, n,$$

con $g(\cdot)$ función dos veces diferenciable y suponga que utilizamos una aproximación de Taylor de primer orden en torno del promedio. Es decir,

$$g(x_i) \approx g(\bar{x}) + g'(\bar{x})(x_i - \bar{x}).$$

De este modo, $\bar{y} = \bar{g}(\mathbf{x})$, y

$$\begin{aligned} \bar{g}(\mathbf{x}) &= \frac{1}{n} \sum_{i=1}^n g(x_i) \approx \frac{1}{n} \sum_{i=1}^n \{g(\bar{x}) + g'(\bar{x})(x_i - \bar{x})\} \\ &= g(\bar{x}) + g'(\bar{x}) \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) = g(\bar{x}). \end{aligned}$$

Mientras que,

$$\begin{aligned} s_y^2 &= \frac{1}{n-1} \sum_{i=1}^n (g(x_i) - \bar{g}(\mathbf{x}))^2 \approx \frac{1}{n-1} \sum_{i=1}^n (g(x_i) - g(\bar{x}))^2 \\ &\approx \frac{1}{n-1} \sum_{i=1}^n (g(\bar{x}) + g'(\bar{x})(x_i - \bar{x}) - g(\bar{x}))^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n (g'(\bar{x})(x_i - \bar{x}))^2 = \{g'(\bar{x})\}^2 \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \{g'(\bar{x})\}^2 s_x^2. \end{aligned}$$

Para $g(x_i) = ax_i + b$ una transformación lineal, tenemos $g'(x_i)$ y luego, se recupera los resultados dados en el ítem (d). Por otro lado, si consideramos una aproximación de Taylor de segundo orden (en torno de \bar{x}),

$$y_i \approx g(\bar{x}) + g'(\bar{x})(x_i - \bar{x}) + \frac{g''(\bar{x})}{2}(x_i - \bar{x})^2.$$

Obtenemos $\bar{y} \approx g(\bar{x}) + g''(\bar{x})s_x^2/2$.

EJEMPLO 1.26. Considere un conjunto de datos x_1, \dots, x_n . Verifique que

$$s^2 = \frac{1}{2n(n-1)} \sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2.$$

Es fácil notar que

$$(x_i - x_j)^2 = ((x_i - \bar{x}) - (x_j - \bar{x}))^2 = (x_i - \bar{x})^2 + (x_j - \bar{x})^2 - 2(x_i - \bar{x})(x_j - \bar{x}).$$

De este modo,

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2 &= \sum_{i=1}^n \sum_{j=1}^n \{(x_i - \bar{x})^2 + (x_j - \bar{x})^2 - 2(x_i - \bar{x})(x_j - \bar{x})\} \\ &= n \sum_{i=1}^n (x_i - \bar{x})^2 + n \sum_{j=1}^n (x_j - \bar{x})^2 - 2 \sum_{i=1}^n (x_i - \bar{x}) \sum_{j=1}^n (x_j - \bar{x}). \end{aligned}$$

Usando la Propiedad 1.25 (a), tenemos $\sum_{i=1}^n (x_i - \bar{x}) = 0$ (y análogamente para $\sum_{j=1}^n (x_j - \bar{x}) = 0$), luego

$$\sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2 = 2n \sum_{i=1}^n (x_i - \bar{x})^2.$$

Lo que lleva a

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{2n(n-1)} \sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2.$$

que es el resultado deseado.

PROPIEDAD 1.27. La mediana es el valor que minimiza la función

$$Q(a) = \sum_{i=1}^n |x_i - a|.$$

DEMOSTRACIÓN. Podemos escribir,

$$Q(a) = \sum_{x_i \geq a} (x_i - a) - \sum_{x_i < a} (x_i - a),$$

así

$$\frac{d}{da} Q(a) = - \sum_{x_i \geq a} 1 + \sum_{x_i < a} 1,$$

resolviendo la condición de primer orden $dQ(a)/da = 0$, tenemos

$$\sum_{x_i \geq a} 1 = \sum_{x_i < a} 1 \quad \implies \quad \hat{a} = \text{me},$$

por la propia definición de mediana. □

OBSERVACIÓN 1.28. Recuerde que

$$z = \text{signo}(z) \cdot |z|.$$

De este modo, $\text{signo}(z) = z/|z|$ (o bien $|z| = \text{signo}(z) \cdot z$). En efecto, tenemos que

$$\frac{d}{dz} |z| = \text{signo}(z), \quad \text{para } z \neq 0.$$

Por tanto, una forma alternativa de escribir $dQ(a)/da$ es

$$\frac{d}{da}Q(a) = \sum_{i=1}^n \text{signo}(x_i - a),$$

es decir

$$\sum_{i=1}^n \text{signo}(x_i - a) = \sum_{i=1}^n \frac{1}{|x_i - a|} (x_i - a) = \sum_{i=1}^n \omega_i(a)(x_i - a), \quad (1.20)$$

donde $\omega_i(a) = |x_i - a|^{-1}$ (Ecuación (1.20) es un promedio ponderado). Sin embargo, resolver la condición de primer orden

$$\sum_{i=1}^n \omega_i(a)(x_i - a) = 0, \quad (1.21)$$

con relación a a es bastante difícil. Una alternativa es considerar el Algoritmo 1 presentado a continuación.

Algoritmo 1: Cálculo de la mediana usando un procedimiento basado en promedios iterativamente ponderados.

Entrada: Conjunto de n datos $\mathbf{x} = (x_1, \dots, x_n)^\top$, una aproximación inicial, a_0 (por ejemplo, $a_0 = \bar{x}$) y un valor de tolerancia ($= 1 \cdot 10^{-6}$)

Salida : Aproximación para el valor de la mediana, $\hat{a} = \text{me}$.

1 **begin**

2 Calcular

$$\omega_i(a_0) \leftarrow \frac{1}{|x_i - a_0|}, \quad i = 1, \dots, n.$$

3 Actualizar:

$$a_1 = \frac{1}{\sum_{j=1}^n \omega_j(a_0)} \sum_{i=1}^n \omega_i(a_0) x_i. \quad (1.22)$$

4 **if** $|a_1 - a_0| < \text{tolerancia}$ **then**

5 **return** $\text{me} = a_1$, y detener el algoritmo.

6 **else**

7 hacer, $a_0 \leftarrow a_1$

8 volver a Paso 2.

9 **end**

10 **end**

Debemos destacar que la etapa de actualización dada en la Ecuación (1.22), equivale a resolver

$$\sum_{i=1}^n \omega_i(a_0)(x_i - a) = 0,$$

con relación a “ a ”, que es muchísimo más simple que resolver (1.21) pues ahora las ponderaciones $\omega_i(a_0)$ están *fijas*.

OBSERVACIÓN 1.29. ¿Qué deficiencias presenta el Algoritmo 1?

- Este algoritmo converge a la mediana (aunque su velocidad de convergencia puede ser bastante lenta).
- En la convergencia, digamos $\hat{a} = \text{me}$, tendremos exactamente un “peso” $\omega_i(\hat{a})$ indefinido (¿Ud. podría ‘adivinar’ cuál?).
- Una alternativa para el punto anterior es usar que¹

$$\text{signo}(z) \approx \frac{z}{\sqrt{z^2 + \epsilon^2}}.$$

Por tanto, “cerca” del óptimo ($\hat{a} = \text{me}$) podemos considerar los *pesos modificados*:

$$\tilde{\omega}_i(a) = \frac{1}{\sqrt{(x_i - a)^2 + \epsilon^2}}, \quad i = 1, \dots, n.$$

EJEMPLO 1.30. Considere el siguiente conjunto de datos:

$$\mathbf{x} = (2.40, 2.70, 2.80, 3.03, 3.40, 3.70, 28.95)^\top.$$

Evidentemente el valor 28.95 es una observación que se destaca y puede ser considerada como atípica. Los valores de la media muestral y el desviación estándar están dados por $\bar{x} = 6.711$ y $s_x = 9.816$, respectivamente. En efecto, es posible especular que el valor de la media muestral no representa una buena estimación del centro de los datos. Suponga que \mathbf{z} denota el conjunto donde hemos eliminado el valor “sospechoso” de 28.95. En este caso obtenemos,

$$\bar{z} = \frac{18.03}{6} = 3.005, \quad s_z^2 = \frac{1.14075}{5} = 0.2282,$$

con $s_z = 0.4777$. Mientras que, cuando calculamos el valor de la mediana para ambos conjuntos de datos obtenemos

$$\text{me}(\mathbf{x}) = 3.03, \quad \text{me}(\mathbf{z}) = \frac{2.80 + 3.03}{2} = 2.915.$$

Es decir, la mediana es un procedimiento muy apropiado para cuantificar el valor central de un conjunto de datos en presencia de observaciones aberrantes o atípicas. Se dice entonces que la mediana es una estadística *robusta*.

Por otro lado, una alternativa robusta a la desviación estándar se conoce como la *desviación mediana absoluta (MAD)*, que es definida como:

$$\text{MAD}(\mathbf{x}) = \text{me}(|\mathbf{x} - \text{me}(\mathbf{x})|).$$

Evaluable esta cantidad para los conjuntos de datos \mathbf{x} y \mathbf{z} , obtenemos

$$\text{MAD}(\mathbf{x}) = \text{me}\{0.63, 0.33, 0.23, 0.00, 0.37, 0.67, 25.92\} = 0.37$$

$$\text{MAD}(\mathbf{z}) = \text{me}\{0.515, 0.215, 0.115, 0.115, 0.485, 0.785\} = \frac{0.215 + 0.485}{2} = 0.35.$$

Para hacer el MAD comparable con la desviación estándar, se define el *MAD normalizado* como:

$$\text{MADN}(\mathbf{x}) = \frac{\text{MAD}(\mathbf{x})}{0.6745}.$$

¹Aproximación que mejora conforme $\epsilon \rightarrow 0$.

De este modo, obtenemos $\text{MADN}(\mathbf{x}) = 0.5486$ mientras que $\text{MADN}(\mathbf{z}) = 0.5189$. Comparativamente con $s_z = 0.4777$, claramente MAD no se ve influenciado fuertemente por la presencia de observaciones atípicas. A continuación presentamos un fragmento de comandos en R para el cálculo de algunas estadísticas de resumen:

```
# introduciendo datos en la consola de R
> x <- c(2.40, 2.70, 2.80, 3.03, 3.40, 3.70, 28.95)
> x
[1] 2.40 2.70 2.80 3.03 3.40 3.70 28.95

# removiendo la 7a observación
> z <- x[-7]
> z
[1] 2.40 2.70 2.80 3.03 3.40 3.70

# cálculo de la media muestral
> mean(x)
[1] 6.711429
> mean(z)
[1] 3.005

# cálculo de la mediana
> median(x)
[1] 3.03
> median(z)
[1] 2.915

# cálculo de la desviación estandar
> sd(x)
[1] 9.815978
> sd(z)
[1] 0.4776505
```

En R está disponible la función `mad`, para el cálculo del MAD o de su versión normalizada.

```
# cálculo del MAD
> abs(x - median(x))
[1] 0.63 0.33 0.23 0.00 0.37 0.67 25.92
> sort(abs(x - median(x)))
[1] 0.00 0.23 0.33 0.37 0.63 0.67 25.92

> mad(x, constant = 1) # MAD
[1] 0.37
> mad(x) # MAD normalizado
[1] 0.548562

> mad(z, constant = 1)
[1] 0.35
> mad(z)
[1] 0.51891
```

DEFINICIÓN 1.31 (Coeficiente de variación). Este coeficiente es una medida que compara la desviación estándar con el promedio de una muestra y es definido como

$$CV = s/\bar{x}, \quad \bar{x} \neq 0.$$

El coeficiente es particularmente útil para comparar dos o más muestras (o grupos). Un valor pequeño para el CV está asociado a una muestra homogénea.

OBSERVACIÓN 1.32. CV es una medida adimensional. Debemos resaltar que, en Econometría, $1/CV$ es conocido como *razón de Sharpe*.

1.2.3. Cálculo del promedio y varianza muestrales. Es interesante notar que una mala implementación computacional puede hacer que un buen algoritmo sea inútil. Un ejemplo de esto son las pobres implementaciones para el cálculo de estadísticas básicas que son ofrecidas en Microsoft Excel. En efecto, [McCullough y Wilson \(1999, 2002, 2005\)](#) reportan una serie de falencias de los procedimientos estadísticos presentes en Excel. Por otro lado, software estadístico como R ([R Core Team, 2019](#)), o bien hojas de cálculo como Gnumeric² disponen de algoritmos confiables (ver, por ejemplo, [Keeling y Pavur, 2007](#)).

La definición de la varianza muestral dada en la Ecuación (1.15) permite sugerir un algoritmo en 2-pasos para el cálculo de s^2 (ver Algoritmo 2). Aunque se ha demostrado que este es un algoritmo estable ([Chan y Lewis, 1979](#); [Chan et al., 1983](#)), puede no ser recomendable para grandes volúmenes de datos debido a que requiere “pasar a través de los datos dos veces”, es decir, requiere usar dos ciclos for. Mientras que la fórmula de Köning en (1.19) lleva a un algoritmo de 1-paso. En efecto, basado en la fórmula:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \sum_{i=1}^n x_i^2 - \frac{n}{n-1} \bar{x}^2,$$

tenemos la implementación dada en el Algoritmo 3. Desafortunadamente, aunque este procedimiento es más veloz que el Algoritmo 2 (2-pasos), es bien sabido que puede llevar a cancelamientos ‘catastróficos’ y por tanto no es un método recomendable (ver [Chan et al., 1983](#); [Barlow, 1993](#)). Para evitar este tipo de dificultades se ha propuesto algoritmos que explotan la definición de la media y varianza muestrales y que solo requieren de pasar por los datos una única vez.

A continuación se describe el *algoritmo online* (1-paso) propuesto por [West \(1979\)](#). Considere una muestra de tamaño n y sea

$$\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i.$$

De este modo, evidentemente tenemos que

$$\begin{aligned} \bar{x}_n &= \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \left(\sum_{i=1}^{n-1} x_i + x_n \right) = \frac{1}{n} \left((n-1)\bar{x}_{n-1} + x_n \right) \\ &= \frac{1}{n} \left(n\bar{x}_{n-1} - \bar{x}_{n-1} + x_n \right) = \bar{x}_{n-1} + \frac{x_n - \bar{x}_{n-1}}{n}. \end{aligned} \quad (1.23)$$

²URL: www.gnumeric.org

Algoritmo 2: Varianza muestral usando un algoritmo de 2-pasos.**Entrada:** Conjunto de n datos $\mathbf{x} = (x_1, \dots, x_n)^\top$.**Salida :** Promedio y varianza muestrales, \bar{x} y s^2 .

```

1 begin
2    $M \leftarrow x_1$ 
3   for  $i = 2$  to  $n$  do
4      $M \leftarrow M + x_i$ 
5   end
6    $M \leftarrow M/n$ 
7    $T \leftarrow (x_1 - M)^2$ 
8   for  $i = 2$  to  $n$  do
9      $T \leftarrow T + (x_i - M)^2$ 
10  end
11   $\bar{x} \leftarrow M$ 
12   $s^2 \leftarrow \frac{1}{n-1}T$ 
13 end

```

Algoritmo 3: Varianza muestral usando un algoritmo de 1-paso.**Entrada:** Conjunto de n datos $\mathbf{x} = (x_1, \dots, x_n)^\top$.**Salida :** Promedio y varianza muestrales, \bar{x} y s^2 .

```

1 begin
2    $M \leftarrow x_1$ 
3    $T \leftarrow x_1^2$ 
4   for  $i = 2$  to  $n$  do
5      $M \leftarrow M + x_i$ 
6      $T \leftarrow T + x_i^2$ 
7   end
8    $\bar{x} \leftarrow M/n$ 
9    $s^2 \leftarrow \frac{1}{n-1}T - \frac{n}{n-1}\bar{x}^2$ 
10 end

```

La base del algoritmo propuesto por [West \(1979\)](#) es la relación recursiva definida en la Ecuación (1.23). Sea $\delta_n = x_n - \bar{x}_{n-1}$, también podemos definir un algoritmo recursivo para el cálculo de la varianza muestral. En efecto, considere

$$\begin{aligned}
T_n &= \sum_{i=1}^n (x_i - \bar{x}_n)^2 = \sum_{i=1}^{n-1} (x_i - \bar{x}_n)^2 + (x_n - \bar{x}_n)^2 \\
&= \sum_{i=1}^{n-1} (x_i - \bar{x}_{n-1} - \delta_n/n)^2 + (x_n - \bar{x}_{n-1} - \delta_n/n)^2 \\
&= \sum_{i=1}^{n-1} \left[(x_i - \bar{x}_{n-1})^2 - 2\frac{\delta_n}{n}(x_i - \bar{x}_{n-1}) + \frac{\delta_n^2}{n^2} + \left(\delta_n - \frac{\delta_n}{n}\right)^2 \right]
\end{aligned}$$

Sabemos por Propiedad 1.25 (a), que

$$\sum_{i=1}^{n-1} (x_i - \bar{x}_{n-1}) = 0,$$

De este modo, podemos escribir la suma de cuadrados T_n , como:

$$\begin{aligned} T_n &= \sum_{i=1}^{n-1} (x_i - \bar{x}_{n-1})^2 + (n-1) \frac{\delta_n^2}{n^2} + \left(1 - \frac{1}{n}\right)^2 \delta_n^2 \\ &= T_{n-1} + \left(1 - \frac{1}{n}\right) \delta_n^2. \end{aligned} \quad (1.24)$$

Ecuaciones (1.23) y (1.24) llevan al Algoritmo 4 (West, 1979), cuya definición se presenta a continuación.

Algoritmo 4: Promedio y varianza muestrales usando un algoritmo on-line.

Entrada: Conjunto de n datos $\mathbf{x} = (x_1, \dots, x_n)^\top$.

Salida : Promedio y varianza muestrales, \bar{x} y s^2 .

```

1 begin
2    $M \leftarrow x_1$ 
3    $T \leftarrow 0$ 
4   for  $i = 2$  to  $n$  do
5      $\delta \leftarrow (x_i - M)/i$ 
6      $M \leftarrow M + \delta$ 
7      $T \leftarrow T + i(i-1)\delta^2$ 
8   end
9    $\bar{x} \leftarrow M$ 
10   $s^2 \leftarrow \frac{1}{n-1} T$ 
11 end
```

Chan y Lewis (1979) introdujeron una medida que permite evaluar la sensibilidad de una muestra $\mathbf{x} = (x_1, \dots, x_n)^\top$ cuando debemos desarrollar el cálculo de su varianza muestral.

Sea

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2},$$

la norma Euclidiana para el conjunto de datos \mathbf{x} . De este modo, se define el *número condición de una muestra*, como:

$$\kappa = \frac{\|\mathbf{x}\|_2}{\sqrt{n-1}s}.$$

Esta medida permite cuantificar el efecto de introducir errores en los datos y como éstos son magnificados en el cálculo de la varianza. Es interesante notar que el número condición puede ser utilizado para evaluar la estabilidad de un algoritmo para el cálculo de s^2 . Además podemos verificar fácilmente que el número condición está relacionado con el coeficiente de variación. Considere el siguiente ejemplo.

EJEMPLO 1.33. Sea $\mathbf{x} = (x_1, \dots, x_n)^\top$ una muestra de datos. Verifique que

$$\kappa^2 = 1 + \frac{n}{n-1} CV^{-2}. \quad (1.25)$$

En efecto, sabemos que

$$\sum_{i=1}^n x_i^2 = \sum_{i=1}^n (x_i - \bar{x})^2 + n\bar{x}^2,$$

dividiendo ámbos términos por $(n-1)s^2$, obtenemos

$$\kappa^2 = \frac{\sum_{i=1}^n x_i^2}{(n-1)s^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n-1)s^2} + \left(\frac{n}{n-1}\right) \frac{\bar{x}^2}{s^2}.$$

Además $\sum_{i=1}^n (x_i - \bar{x})^2 = (n-1)s^2$, entonces podemos escribir

$$\kappa^2 = 1 + \left(\frac{n}{n-1}\right) \left(\frac{\bar{x}}{s}\right)^2.$$

tenemos $CV = s/\bar{x}$ y el resultado sigue.

Desde (1.25), sigue que

$$\kappa = \sqrt{1 + n CV^{-2} / (n-1)},$$

es decir, a menos que CV sea bastante grande $\kappa \approx CV^{-1}$ será una buena aproximación.

Finalmente debemos destacar que Chan et al. (1983) propusieron un algoritmo extremadamente estable para el cálculo de la varianza basado en un procedimiento de *suma acumulada por pares* el cual puede ser fácilmente paralelizado.

1.2.4. Medidas de forma. Momentos de orden mayor, o estadísticas involucrando potencias de orden mayor de los datos observados permiten caracterizar la forma de la densidad que describe el mecanismo que genera las n realizaciones x_1, \dots, x_n de nuestra variable de interés. A continuación revisamos la definición de los coeficientes de asimetría y curtosis, los que caracterizan el grado de asimetría de una distribución en torno de su promedio y el grado de agudeza o achatamiento de una distribución al ser comparada contra la distribución normal (gaussiana).

DEFINICIÓN 1.34 (Coeficiente de asimetría). Considere M_3 el tercer momento muestral en torno del promedio. Entonces, se define el *coeficiente de asimetría* (o sesgo) como:

$$b_1 = \frac{M_3}{s^3} = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^3.$$

OBSERVACIÓN 1.35. b_1 también es conocido como coeficiente de asimetría de Fisher y se caracteriza por ser una medida adimensional así como por ser invariante bajo traslaciones del origen y transformaciones de escala. Considere los gráficos desplegados en la Figura 1

- Si $b_1 = 0$ la distribución es simétrica con relación a la media.
- Si $b_1 > 0$ la distribución tiene sesgo positivo (o hacia la derecha), en cuyo caso la distribución tiende a concentrarse en valores altos de la variable. En caso contrario ($b_1 < 0$), diremos que su sesgo es negativo (o que es asimétrica hacia la izquierda).

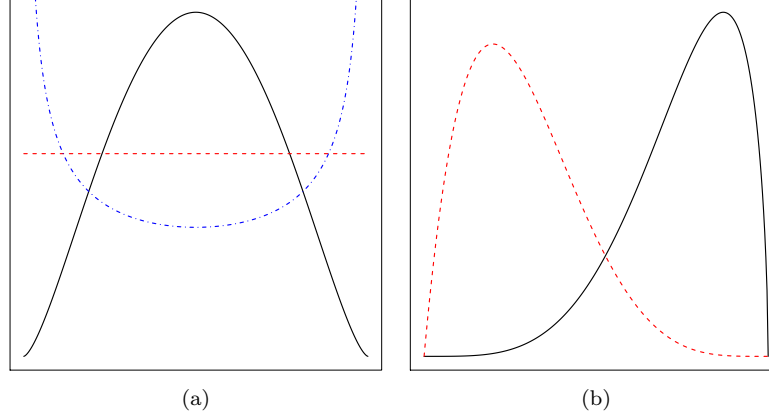


Figura 1. (a) distribuciones simétricas, (b) distribución con asimetría negativa (—) y asimetría positiva (— —).

OBSERVACIÓN 1.36. Se han definido varios índices de simetría, por ejemplo el *coeficiente de asimetría de Galton*:

$$b_G = \frac{(Q_3 - Q_2) - (Q_2 - Q_1)}{Q_3 - Q_1}.$$

DEFINICIÓN 1.37 (Coeficiente de curtosis). Considere M_4 el cuarto momento muestral en torno del promedio. Entonces, se define el *coeficiente de curtosis* (o achataamiento) como:

$$b_2 = \frac{M_4}{s^4} - 3 = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^4 - 3.$$

OBSERVACIÓN 1.38. El término -3 hace que $b_2 = 0$ cuando los datos siguen una distribución normal (gaussiana) en cuyo caso decimos que la distribución es *mesocúrtica*. Considere la Figura 2, si $b_2 > 0$ la distribución de los datos es más aguzada que la distribución normal (distribución *leptocúrtica*), mientras que si $b_2 < 0$ la distribución es más achatada que la normal (distribución *platicúrtica*).

OBSERVACIÓN 1.39. Debemos destacar que Spicer (1972) propuso calcular momentos centrales de hasta cuarto orden, M_2, M_3 y M_4 usando un algoritmo online. La función `moments` desde la biblioteca `fastmatrix` permite el cálculo de los coeficientes b_1 y b_2 .

1.2.5. Estadísticas de resumen para datos agrupados. Cuando los datos han sido organizados en una *tabla de frecuencias* se dice que los datos se encuentran agrupados. El objetivo de esta sección es proporcionar fórmulas para las medidas de posición (o tendencia central), de dispersión y de forma sin la necesidad de desagregar los datos. Primeramente vamos a suponer que la variable de interés es discreta. Considere el siguiente ejemplo:

EJEMPLO 1.40. Se consultó las fichas de los empleados de una fábrica, registrándose el *número de cargas familiares* y se obtuvo los siguientes datos:

1 2 4 2 2 2 3 2 1 1 0 2 2
0 2 2 1 2 2 3 1 2 2 1 2

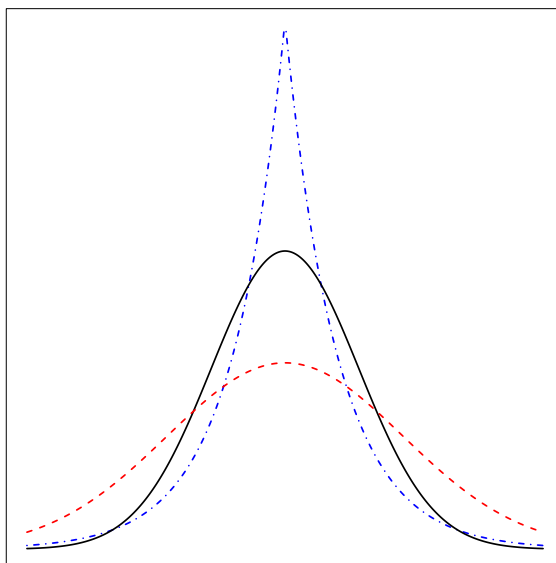


Figura 2. Distintos grados de curtosis: distribución leptocúrtica (---), mesocúrtica (—) y platicúrtica (---).

En cuyo caso, tenemos que la tabla de frecuencias asociada asume la forma:

Cargas familiares	Número de empleados	Porcentaje de empleados	Num. acumulado de empleados	Porcentaje acumulado
0	2	8 %	2	8 %
1	6	23 %	8	32 %
2	14	56 %	22	88 %
3	2	8 %	24	96 %
4	1	4 %	25	100 %
Total	25	100 %	—	—

Podemos notar que los datos de interés corresponden a una variable discreta x (número de cargas familiares) que tiene k categorías. De este modo, podemos construir la siguiente tabla

Variable	Frecuencia Absoluta	Frecuencia Relativa	Frec. Abs. Acumulada	Frec. Rel. Acumulada
x_1	n_1	f_1	N_1	F_1
x_2	n_2	f_2	N_2	F_2
\vdots	\vdots	\vdots	\vdots	\vdots
x_k	n_k	f_k	N_k	F_k
Total	n	1	—	—

donde

$$f_i = \frac{n_i}{n}, \quad N_i = \sum_{j=1}^i n_j, \quad F_i = \sum_{j=1}^i f_j, \quad (1.26)$$

para $i = 1, \dots, k$. Evidentemente tenemos que

$$\sum_{i=1}^k n_i = n, \quad N_k = n, \quad \sum_{i=1}^k f_i = 1, \quad F_k = 1. \quad (1.27)$$

En este contexto tenemos que la media y varianza muestrales están dados por

$$\bar{x} = \frac{1}{n} \sum_{i=1}^k n_i x_i = \sum_{i=1}^k f_i x_i,$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^k n_i (x_i - \bar{x})^2.$$

Para los datos del Ejemplo 1.40, es fácil notar que

$$\bar{x} = \frac{2 \cdot 0 + 6 \cdot 1 + 14 \cdot 2 + 2 \cdot 3 + 1 \cdot 4}{25} = \frac{44}{25} = 1.760$$

$$s^2 = \frac{2(0 - 1.76)^2 + 6(1 - 1.76)^2 + 14(2 - 1.76)^2 + 2(3 - 1.76)^2 + 1(4 - 1.76)^2}{25 - 1}$$

$$= \frac{18.56}{24} = 0.773$$

Para datos continuos x_1, \dots, x_n tenemos la tabla de frecuencias:

Marca de clase	Frecuencia Absoluta	Frecuencia Relativa	Frec. Abs. Acumulada	Frec. Rel. Acumulada
C_1	n_1	f_1	N_1	F_1
C_2	n_2	f_2	N_2	F_2
\vdots	\vdots	\vdots	\vdots	\vdots
C_k	n_k	f_k	N_k	F_k
Total	n	1	—	—

donde f_i , N_i y F_i , para cada una de las k categorías ha sido definido en (1.26) y (1.27), mientras que la *marca de clase* es definida como:

$$C_i = \frac{L_i + U_i}{2}, \quad i = 1, \dots, k,$$

con L_i y U_i los límites inferior y superior de cada intervalo, respectivamente. Note que la marca de clase es un representante de la clase (intervalo) respectiva. Análogamente a caso de una tabla de frecuencias para datos discretos, tenemos que:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^k n_i C_i = \sum_{i=1}^k f_i C_i,$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^k n_i (C_i - \bar{x})^2.$$

EJEMPLO 1.41. Considere los *precios de cierre* de acciones una determinada empresa nacional, dados por:

179	173	181	170	158	174	172	166	194	185
162	187	198	177	178	165	154	188	166	171
175	182	167	169	172	186	172	176	168	187

cuya tabla de frecuencias asume la forma:

Precio de cierres	Marca de clase	Número de días	Porcentaje de días	Num. de días acumulado	Porcentaje acumulado
(150,160]	155	2	7 %	2	7 %
(160,170]	165	8	27 %	10	34 %
(170,180]	175	11	36 %	21	70 %
(180,190]	185	7	23 %	28	93 %
(190,200]	195	2	7 %	30	100 %
Total	—	30	100 %	—	—

De este modo, podemos calcular la media y varianza muestrales desde la tabla de frecuencias, como:

$$\bar{x} = \frac{2 \cdot 155 + 8 \cdot 165 + 11 \cdot 175 + 7 \cdot 185 + 2 \cdot 195}{30} = \frac{5240}{30} = 174.667$$

mientras que,

$$\begin{aligned} \sum_{i=1}^5 n_i (C_i - \bar{x})^2 &= 2(155 - 174.667)^2 + 8(165 - 174.667)^2 + 11(175 - 174.667)^2 \\ &\quad + 7(185 - 174.667)^2 + 2(195 - 174.667)^2 = 3096.667, \end{aligned}$$

de este modo, $s^2 = 3096.667/(30-1) = 106.782$. Es interesante notar que los valores de la media y varianza muestrales obtenidos a partir de una tabla de frecuencias corresponden a una *aproximación* de los valores obtenidos a partir de los *datos a granel*. En efecto, el siguiente fragmento en R considera los datos crudos:

```
# precios de cierre (datos a granel)
> x <- c(179, 173, 181, 170, 158, 174, 172, 166, 194, 185,
+       162, 187, 198, 177, 178, 165, 154, 188, 166, 171, 175,
+       182, 167, 169, 172, 186, 172, 176, 168, 187)

# cálculo de media y varianza muestrales
> mean(x)
[1] 175.0667
> var(x)
[1] 105.7195
```

Se ha sugerido una serie de procedimientos para la construcción de tablas de frecuencia. Considere los siguientes pasos:

- Determinar el *número de categorías*, k usando por ejemplo:

$$k = \sqrt{n},$$

$$k = 1 + 3.3 \log_{10}(n), \quad (\text{regla de Sturges})$$

$$k = \lceil 2n^{1/3} \rceil, \quad (\text{regla de Rice})$$

- Calcular el rango,

$$R = \max\{x_i\}_{i=1}^n - \min\{x_i\}_{i=1}^n,$$

- Determinar la *longitud de los intervalos* a_i para $i = 1, \dots, k$. Usualmente, podemos elegir una longitud constante, $a_i = a$ como:

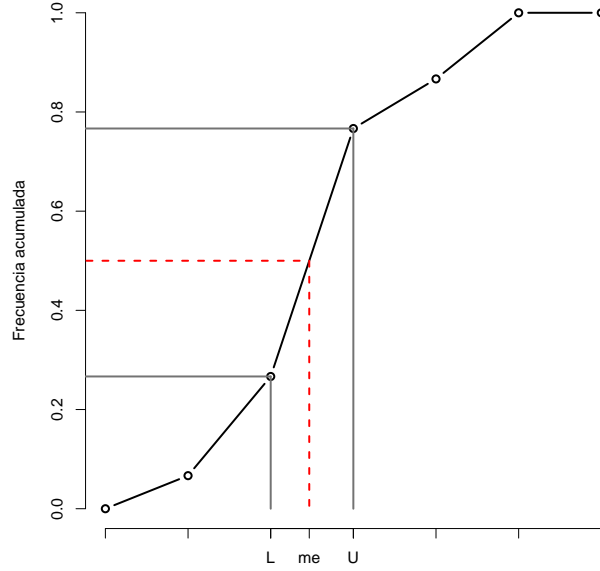
$$a = \frac{R}{k},$$

- Finalmente, construir los *límites de clases*

$$\begin{aligned} L_1 &= \min\{x_i\}_{i=1}^n - \Delta, & U_1 &= L_1 + a, \\ L_2 &= U_1, & U_2 &= L_2 + a, \\ &\vdots & & \end{aligned}$$

Debemos resaltar que aunque este tipo de procedimientos suelen ser bastante apropiados en la construcción de tablas de frecuencia, deben ser entendidos meramente como *reglas de trabajo*.

Considere la función de distribución acumulada



de este modo, podemos usar interpolación lineal para determinar la mediana en datos agrupados. En efecto,

$$\frac{1/2 - F_{i-1}}{F_i - F_{i-1}} = \frac{\text{me} - L_i}{U_i - L_i},$$

es decir,

$$(1/2 - F_{i-1})(U_i - L_i) = (\text{me} - L_i)(F_i - F_{i-1}),$$

o bien

$$\text{me} = L_i + \frac{1/2 - F_{i-1}}{F_i - F_{i-1}}(U_i - L_i) = L_i + \frac{1/2 - F_{i-1}}{F_i - F_{i-1}} a_i,$$

donde $a_i = U_i - L_i$, representa la amplitud de la clase mediana. Recordando que $F_i = \sum_{j=1}^i f_j$, tenemos

$$F_i - F_{i-1} = \sum_{j=1}^i f_j - \sum_{j=1}^{i-1} f_j = f_i + \sum_{j=1}^{i-1} f_j - \sum_{j=1}^{i-1} f_j = f_i.$$

Además, como $f_i = n_i/n$ y $F_i = N_i/n$, podemos re-escribir la *mediana para datos agrupados* como:

$$\begin{aligned} \text{me} &= L_i + \frac{1/2 - F_{i-1}}{f_i} a_i = L_i + \frac{1/2 - F_{i-1}}{n_i/n} a_i \\ &= L_i + \frac{n/2 - N_{i-1}}{n_i} a_i \end{aligned}$$

Siguiendo exactamente la misma lógica, tenemos que el percentil k -ésimo, digamos P_k es dado por

$$P_k = L_i + \frac{k/100 - F_{i-1}}{f_i} a_i,$$

o bien

$$P_k = L_i + \frac{n(k/100) - N_{i-1}}{n_i} a_i.$$

EJEMPLO 1.42. Los trabajadores de una empresa, cuya tarea es clasificar y envasar fruta, obtuvieron los siguientes salarios semanales (clasificados según sexo).

Ingreso (UM)	Mujeres	Hombres
65 – 75	10	0
75 – 85	15	0
85 – 95	60	5
95 – 105	15	10
105 – 115	10	50
115 – 125	0	25
125 – 135	0	10
Total	110	100

Consideraremos solamente el grupo de mujeres y dejaremos el análisis del grupo de hombres y el total de trabajadores como ejercicio. Así, la tabla de frecuencias para el grupo de mujeres adopta la forma:

Ingreso (UM)	C_i	n_i	f_i	N_i	F_i
65 – 75	70	10	0.090	10	0.090
75 – 85	80	15	0.136	25	0.226
85 – 95	90	60	0.548	85	0.774
95 – 105	100	15	0.136	100	0.910
105 – 115	110	10	0.090	110	1.000
Total	–	110	1.000	–	–

En este caso, tenemos $n = 110$ y $\sum_{i=1}^5 n_i C_i = 9900$. De este modo, la media aritmética para el grupo de mujeres es:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^5 n_i C_i = \frac{9900}{110} = 90 \text{ (UM)}.$$

Para calcular la mediana, primero debemos ubicar el intervalo mediano. En efecto, debemos ubicar la primera frecuencia relativa acumulada (F_i) que supere 0.5 (o

bien, frecuencia absoluta acumulada (N_i) que supere $n/2$). De este modo el intervalo mediano es $(85, 95]$. Además, $a_i = 10$ para todos los intervalos. De este modo,

$$\text{me} = L_i + \frac{1/2 - F_{i-1}}{f_i} a_i,$$

donde $L_i = 85$, $F_{i-1} = 0.226$, $f_i = 0.548$ y $a_i = 10$, luego

$$\text{me} = 85 + \frac{0.500 - 0.226}{0.548} \cdot 10 = 85 + 0.5 \cdot 10 = 90 \text{ (UM)}.$$

Por otro lado,

$$s^2 = \frac{1}{n-1} \left(\sum_{i=1}^5 n_i C_i^2 - n \bar{x}^2 \right).$$

En nuestro caso,

$$\sum_{i=1}^5 n_i C_i^2 = 902\,000, \quad \bar{x}^2 = 8\,100.$$

Así,

$$\begin{aligned} s^2 &= \frac{1}{110-1} (902\,000 - 110 \cdot 8\,100^2) = \frac{1}{109} (902\,000 - 891\,000) \\ &= \frac{11\,000}{109} = 100.9174 \text{ (UM)}^2. \end{aligned}$$

Además, tenemos que $s = \sqrt{11\,000/109} = 10.0458$ (UM). Mientras que

$$\text{CV} = \frac{10.0458}{90} = 0.1116.$$

Podemos evaluar la simetría usando el coeficiente de Galton. Por tanto, debemos calcular Q_1 y Q_3 , como:

$$\begin{aligned} Q_1 &= 85 + \frac{0.250 - 0.226}{0.548} \cdot 10 = 85.438 \\ Q_3 &= 85 + \frac{0.750 - 0.226}{0.548} \cdot 10 = 94.562, \end{aligned}$$

de este modo $IQR = 9.1240$, y

$$\begin{aligned} \gamma_G &= \frac{(Q_3 - Q_2) - (Q_2 - Q_1)}{Q_3 - Q_1} \\ &= \frac{(94.562 - 90) - (90 - 85.438)}{9.124} = \frac{4.562 - 4.562}{9.124} = 0.000 \end{aligned}$$

Es decir, la distribución de los datos es simétrica.

Mientras que la moda interpolada, es dada por

$$\text{mo} = L_m + \frac{\Delta_1}{\Delta_1 + \Delta_2} a_m,$$

donde

- L_m es el límite inferior de la clase modal.
- $\Delta_1 = n_m - n_{m-1}$ con n_m la frecuencia absoluta de la clase modal, mientras que n_{m-1} es la frecuencia absoluta de la clase anterior a la clase modal.
- $\Delta_2 = n_m - n_{m+1}$ con n_{m+1} la frecuencia absoluta de la clase posterior a la clase modal.
- a_m es la amplitud de la clase modal.

OBSERVACIÓN 1.43. La clase modal es aquella que tiene la mayor frecuencia relativa. Note además, que en una tabla de frecuencia podría existir más de una clase modal.

1.3. Covarianza y correlación

Considere $\mathbf{x} = (x_1, \dots, x_n)^\top$ y $\mathbf{y} = (y_1, \dots, y_n)^\top$ dos vectores de datos con n observaciones. En esta sección se introducen medidas de asociación entre variables continuas y ordinales. Primeramente se introduce el concepto de covarianza.

DEFINICIÓN 1.44 (Covarianza). Para el conjunto $(x_1, y_1), \dots, (x_n, y_n)$, se define la covarianza como una medida de variabilidad conjunta de dos variables cuantitativas, como:

$$\text{cov}(\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

OBSERVACIÓN 1.45. En efecto, es fácil apreciar que $\text{cov}(\mathbf{x}, \mathbf{x}) = s_x^2$. Además, evidentemente la covarianza está relacionada con el producto interno entre dos vectores n -dimensionales, $\langle \mathbf{u}, \mathbf{v} \rangle = \sum_{i=1}^n u_i v_i$.

EJEMPLO 1.46. Tal como en la Propiedad 1.25 (b), a continuación derivamos una expresión alternativa para la covarianza. En efecto,

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) &= \sum_{i=1}^n x_i(y_i - \bar{y}) - \bar{x} \sum_{i=1}^n (y_i - \bar{y}) \\ &= \sum_{i=1}^n x_i y_i - \bar{y} \sum_{i=1}^n x_i = \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}. \end{aligned}$$

De este modo, podemos escribir

$$\text{cov}(\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \left(\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} \right).$$

PROPIEDAD 1.47. Tenemos que

$$\text{cov}(a\mathbf{x} + b, c\mathbf{y} + d) = ac \text{cov}(\mathbf{x}, \mathbf{y}).$$

DEMOSTRACIÓN. Sea $z_i = ax_i + b$ y $w_i = cy_i + d$ para $i = 1, \dots, n$. Entonces $\bar{z} = a\bar{x} + b$ y $\bar{w} = c\bar{y} + d$, luego

$$\text{cov}(\mathbf{z}, \mathbf{w}) = \frac{1}{n-1} \sum_{i=1}^n (z_i - \bar{z})(w_i - \bar{w}),$$

como $z_i - \bar{z} = a(x_i - \bar{x})$ y $w_i - \bar{w} = c(y_i - \bar{y})$. Entonces

$$\text{cov}(\mathbf{z}, \mathbf{w}) = \frac{1}{n-1} ac \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}),$$

que es el resultado deseado. \square

La covarianza permite medir la asociación, pero depende de la unidad de medida. Una alternativa es usar una medida conocida como *correlación*.

DEFINICIÓN 1.48 (Correlación). La correlación entre \mathbf{x} e \mathbf{y} es la covarianza de sus versiones estandarizadas. Es decir,

$$\begin{aligned}\text{cor}(\mathbf{x}, \mathbf{y}) &= \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right) \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{s_x s_y}.\end{aligned}$$

OBSERVACIÓN 1.49. $\text{cor}(\mathbf{x}, \mathbf{y})$ es una medida adimensional.

EJERCICIO 1.50. Muestre que

$$\text{cor}(a\mathbf{x} + b, c\mathbf{y} + d) = \pm \text{cor}(\mathbf{x}, \mathbf{y}).$$

EJEMPLO 1.51. Verifique que

$$s_{x+y}^2 = s_x^2 + s_y^2 + 2 \text{cov}(\mathbf{x}, \mathbf{y}) = s_x^2 + s_y^2 + 2 \text{cor}(\mathbf{x}, \mathbf{y}) s_x s_y.$$

En efecto, podemos notar que

$$\begin{aligned}s_{x+y}^2 &= \frac{1}{n-1} \sum_{i=1}^n ((x_i - \bar{x}) + (y_i - \bar{y}))^2 \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^n (y_i - \bar{y})^2 + 2 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right) \\ &= s_x^2 + s_y^2 + 2 \text{cov}(\mathbf{x}, \mathbf{y})\end{aligned}$$

OBSERVACIÓN 1.52. Cuando $\text{cor}(\mathbf{x}, \mathbf{y}) = 0$ diremos que \mathbf{x} e \mathbf{y} son *no correlacionados*.

OBSERVACIÓN 1.53. Si \mathbf{x} e \mathbf{y} son no correlacionados, entonces

$$s_{x+y}^2 = s_x^2 + s_y^2.$$

PROPIEDAD 1.54. Tenemos que $\{\text{cor}(\mathbf{x}, \mathbf{y})\}^2 \leq 1$.

DEMOSTRACIÓN. Usando la desigualdad de Cauchy-Schwarz³, tenemos que

$$\left(\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right)^2 \leq \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right) \left(\sum_{i=1}^n (y_i - \bar{y})^2 \right),$$

es decir,

$$\begin{aligned}\left(\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right)^2 &\leq \left(\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right) \left(\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 \right) \\ \{\text{cov}(\mathbf{x}, \mathbf{y})\}^2 &\leq s_x^2 s_y^2.\end{aligned}$$

De este modo,

$$\frac{\text{cov}^2(\mathbf{x}, \mathbf{y})}{s_x^2 s_y^2} \leq 1.$$

□

³Recuerde que $\langle \mathbf{u}, \mathbf{v} \rangle^2 \leq \|\mathbf{u}\|^2 \|\mathbf{v}\|^2$

Evidentemente, desde la propiedad anterior tenemos $-1 \leq \text{cor}(\mathbf{x}, \mathbf{y}) \leq 1$.

Para datos que no son de naturaleza continua, una alternativa es reemplazar el valor de x_i (y_i) por su rango, R_i (S_i) que corresponde a un valor en el conjunto $\{1, 2, \dots, n\}$. Considere la siguiente definición

DEFINICIÓN 1.55 (Coeficiente de correlación de Spearman). Suponga el conjunto de datos pareados $(x_1, y_1), \dots, (x_n, y_n)$ y sea R_i el rango de x_i y S_i el rango de y_i , para $i = 1, \dots, n$. Entonces el coeficiente de correlación de Spearman es dado por

$$r_S = \frac{\sum_{i=1}^n (R_i - \bar{R})(S_i - \bar{S})}{\sqrt{\sum_{i=1}^n (R_i - \bar{R})^2 \sum_{i=1}^n (S_i - \bar{S})^2}}.$$

Considere

$$D = \sum_{i=1}^n (R_i - S_i)^2,$$

y suponga que no existen empates entre los x 's e y 's, entonces podemos escribir

$$r_S = 1 - \frac{6D}{n(n^2 - 1)}.$$

OBSERVACIÓN 1.56. Análogamente a los resultados expuestos en la Sección 1.2.3, podemos llevar a cabo el cálculo de la covarianza entre \mathbf{x} e \mathbf{y} de manera eficiente considerando,

$$C_n = \sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n) = \sum_{i=1}^{n-1} (x_i - \bar{x}_n)(y_i - \bar{y}_n) + (x_n - \bar{x}_n)(y_n - \bar{y}_n).$$

Sea $\delta_n = x_n - \bar{x}_{n-1}$, de este modo podemos escribir:

$$x_i - \bar{x}_n = x_i - \bar{x}_{n-1} - \frac{\delta_n}{n}, \quad x_n - \bar{x}_n = \left(1 - \frac{1}{n}\right)\delta_n,$$

y análogamente,

$$y_i - \bar{y}_n = y_i - \bar{y}_{n-1} - \frac{\eta_n}{n}, \quad y_n - \bar{y}_n = \left(1 - \frac{1}{n}\right)\eta_n,$$

con $\eta_n = y_n - \bar{y}_{n-1}$. Notando que

$$\begin{aligned} (x_i - \bar{x}_n)(y_i - \bar{y}_n) &= \left\{ (x_i - \bar{x}_{n-1}) - \frac{\delta_n}{n} \right\} \left\{ (y_i - \bar{y}_{n-1}) - \frac{\eta_n}{n} \right\} \\ &= (x_i - \bar{x}_{n-1})(y_i - \bar{y}_{n-1}) - (x_i - \bar{x}_{n-1})\frac{\eta_n}{n} \\ &\quad - \frac{\delta_n}{n}(y_i - \bar{y}_{n-1}) + \frac{\delta_n \eta_n}{n^2}. \end{aligned}$$

Sumando sobre $\{1, \dots, n-1\}$ y recordando que

$$\sum_{i=1}^{n-1} (x_i - \bar{x}_{n-1}) = 0, \quad \sum_{i=1}^{n-1} (y_i - \bar{y}_{n-1}) = 0,$$

sigue que

$$\begin{aligned} C_n &= \sum_{i=1}^{n-1} (x_i - \bar{x}_{n-1})(y_i - \bar{y}_{n-1}) + \left(\frac{n-1}{n^2}\right)\delta_n \eta_n + \left(\frac{n-1}{n}\right)^2 \delta_n \eta_n \\ &= C_{n-1} + \left(\frac{n-1}{n}\right)\delta_n \eta_n. \end{aligned}$$

Así, haciendo $\text{cov}_n = \text{cov}(\mathbf{x}, \mathbf{y})$, tenemos

$$\text{cov}_n = \frac{1}{n-1} C_n = \frac{1}{n-1} \left\{ C_{n-1} + \left(\frac{n-1}{n} \right) \delta_n \eta_n \right\},$$

pero $C_{n-1} = ((n-1) - 1) \text{cov}_{n-1} = (n-1) \text{cov}_{n-1} - \text{cov}_{n-1}$. De este modo,

$$\begin{aligned} \text{cov}_n &= \frac{1}{n-1} \left\{ (n-1) \text{cov}_{n-1} - \text{cov}_{n-1} + \left(\frac{n-1}{n} \right) \delta_n \eta_n \right\} \\ &= \text{cov}_{n-1} - \frac{\text{cov}_{n-1}}{n-1} + \frac{(x_n - \bar{x}_{n-1})(y_n - \bar{y}_{n-1})}{n}. \end{aligned}$$

Este desarrollo lleva al siguiente algoritmo:

Algoritmo 5: Covarianza muestral usando un algoritmo online.

Entrada: Conjuntos de n datos $\mathbf{x} = (x_1, \dots, x_n)^\top$ e $\mathbf{y} = (y_1, \dots, y_n)^\top$.
Salida : Covarianza muestral, $\text{cov}(\mathbf{x}, \mathbf{y})$.

```

1 begin
2    $M \leftarrow x_1$ 
3    $N \leftarrow y_1$ 
4    $C \leftarrow 0$ 
5   for  $i = 2$  to  $n$  do
6      $\delta \leftarrow (x_i - M)/i$ 
7      $\eta \leftarrow (y_i - N)/i$ 
8      $M \leftarrow M + \delta$ 
9      $N \leftarrow N + \eta$ 
10     $C \leftarrow C - \frac{C}{i-1} + i \delta \eta$ 
11  end
12   $\text{cov}(\mathbf{x}, \mathbf{y}) \leftarrow C$ .
13 end
```

Además, es fácil notar que modificando ligeramente el Algoritmo 5 podemos obtener también \bar{x}_n , \bar{y}_n , s_x^2 y s_y^2 .

1.4. Estadísticas descriptivas multivariadas

Considere que nuestro interés es estudiar $p \geq 2$ variables (características) de interés asociadas a una muestra aleatoria $\mathbf{x}_1, \dots, \mathbf{x}_n$ donde cada $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$ es un vector p -dimensional. Note que, podemos disponer la información en una *matriz de datos*

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix}.$$

Análogamente a la media y varianza muestrales \bar{x} y s^2 unidimensionales, respectivamente, podemos definir sus contrapartes multivariadas como:

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i,$$

$$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top.$$

Además, la matriz de correlación entre las p variables de interés, es dada por:

$$\mathbf{R} = (r_{ij}),$$

donde

$$r_{ij} = \frac{\sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j)}{\sqrt{\sum_{k=1}^n (x_{ki} - \bar{x}_i)^2} \sqrt{\sum_{k=1}^n (x_{kj} - \bar{x}_j)^2}} = \frac{s_{ij}}{\sqrt{s_{ii}s_{jj}}},$$

con $\mathbf{S} = (s_{ij})$ y $\bar{x}_i = (\sum_{k=1}^n x_{ki})/n$. Sea, $\mathbf{D} = \text{diag}(s_{11}, s_{22}, \dots, s_{pp})$, así podemos escribir

$$\mathbf{R} = \mathbf{D}^{-1/2} \mathbf{S} \mathbf{D}^{-1/2}.$$

OBSERVACIÓN 1.57. Algunas propiedades del vector de medias y la matriz de covarianza, surgen de escribir formas compactas que dependen de la matriz de datos $\mathbf{X} = (x_{ij}) \in \mathbb{R}^{n \times p}$. En efecto,

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \frac{1}{n} \mathbf{X}^\top \mathbf{1}.$$

Además,

$$\begin{aligned} \mathbf{Q} &= \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top = \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) \mathbf{x}_i^\top - \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) \bar{\mathbf{x}}^\top \\ &= \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - \bar{\mathbf{x}} \sum_{i=1}^n \mathbf{x}_i^\top = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - n \bar{\mathbf{x}} \bar{\mathbf{x}}^\top \end{aligned}$$

Es fácil notar que $\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{X}$, de este modo,

$$\begin{aligned} \mathbf{S} &= \frac{1}{n-1} \mathbf{Q} = \frac{1}{n-1} \left\{ \mathbf{X}^\top \mathbf{X} - n \left(\frac{1}{n} \mathbf{X}^\top \mathbf{1} \right) \left(\frac{1}{n} \mathbf{X}^\top \mathbf{1} \right)^\top \right\} \\ &= \frac{1}{n-1} \left(\mathbf{X}^\top \mathbf{X} - \frac{1}{n} \mathbf{X}^\top \mathbf{1} \mathbf{1}^\top \mathbf{X} \right) = \frac{1}{n-1} \mathbf{X}^\top \mathbf{C} \mathbf{X} \end{aligned}$$

con $\mathbf{C} = \mathbf{I}_n - \frac{1}{n} \mathbf{1} \mathbf{1}^\top$ la matriz de centrado. Es sencillo mostrar que

$$\mathbf{C}^\top = \mathbf{C}, \quad \mathbf{C}^2 = \mathbf{C},$$

es decir \mathbf{C} es matriz de proyección. Esto lleva al siguiente resultado.

PROPIEDAD 1.58. La matriz de covarianza \mathbf{S} , es semidefinida positiva.

DEMOSTRACIÓN. Sea $\mathbf{a} \in \mathbb{R}^p$, vector no nulo. Tenemos que,

$$\begin{aligned} \mathbf{a}^\top \mathbf{S} \mathbf{a} &= \frac{1}{n-1} \mathbf{a}^\top \mathbf{X}^\top \mathbf{C} \mathbf{X} \mathbf{a} = \frac{1}{n-1} \mathbf{a}^\top \mathbf{X}^\top \mathbf{C}^2 \mathbf{X} \mathbf{a} \\ &= \frac{1}{n-1} \mathbf{u}^\top \mathbf{u} \geq 0, \quad \mathbf{u} = \mathbf{C} \mathbf{X} \mathbf{a}, \end{aligned}$$

es decir, \mathbf{S} es matriz semidefinida positiva.⁴ □

Considere la siguiente transformación:

$$\mathbf{y}_i = \mathbf{A}\mathbf{x}_i + \mathbf{b}, \quad i = 1, \dots, n.$$

Entonces, $\bar{\mathbf{y}} = \mathbf{A}\bar{\mathbf{x}} + \mathbf{b}$, mientras que

$$\mathbf{y}_i - \bar{\mathbf{y}} = \mathbf{A}\mathbf{x}_i + \mathbf{b} - \mathbf{A}\bar{\mathbf{x}} - \mathbf{b} = \mathbf{A}(\mathbf{x}_i - \bar{\mathbf{x}}).$$

De este modo,

$$\begin{aligned} \mathbf{S}_y &= \frac{1}{n-1} \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})^\top = \frac{1}{n-1} \sum_{i=1}^n \mathbf{A}(\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{A}^\top \\ &= \frac{1}{n-1} \mathbf{A} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{A}^\top = \mathbf{A}\mathbf{S}_x\mathbf{A}^\top. \end{aligned}$$

En particular, para la transformación (de Mahalanobis),

$$\mathbf{z}_i = \mathbf{S}^{-1/2}(\mathbf{x}_i - \bar{\mathbf{x}}), \quad i = 1, \dots, n,$$

donde $\mathbf{S} = \mathbf{S}^{1/2}\mathbf{S}^{1/2}$ con $\mathbf{S}^{1/2}$ un factor raíz cuadrada de \mathbf{S} . Entonces, sigue que

$$\bar{\mathbf{z}} = \mathbf{0}, \quad \text{y} \quad \mathbf{S}_z = \mathbf{I}_p.$$

DEFINICIÓN 1.59 (Distancia de Mahalanobis). Considere una muestra de n observaciones $\mathbf{x}_1, \dots, \mathbf{x}_n$ de este modo, la distancia de Mahalanobis es dada por

$$D_i = \{(\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{S}^{-1}(\mathbf{x}_i - \bar{\mathbf{x}})\}^{1/2}, \quad i = 1, \dots, n,$$

como la distancia de la observación i -ésima hacia el “centro” de los datos, $\bar{\mathbf{x}}$ ponderada por la matriz de covarianza.

Sea

$$g_{ij} = (\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{S}^{-1}(\mathbf{x}_j - \bar{\mathbf{x}}), \quad i, j = 1, \dots, n.$$

lo que permite definir medidas de sesgo y curtosis multivariadas (Mardia, 1970), dadas por

$$b_{1p} = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n g_{ij}^3, \quad b_{2p} = \frac{1}{n} \sum_{i=1}^n g_{ii}^2,$$

respectivamente. Las cantidades anteriores llevan a un test para evaluar la normalidad multivariada (ver Mardia, 1974).

EJEMPLO 1.60. El conjunto de datos de flores Iris (Fisher, 1936), corresponden a 150 mediciones (en centímetros) del largo y ancho de los sépalos y largo y ancho de los pétalos para flores iris de las especies Setosa, Versicolor y Virginica. Los siguientes comandos en R permiten obtener el diagrama de dispersión así como algunas estadísticas de resumen para los datos de Iris

```
# datos Iris (50 observaciones por cada especie)
> iris
  Sepal.Length Sepal.Width Petal.Length Petal.Width Species
1         5.1         3.5         1.4         0.2    setosa
2         4.9         3.0         1.4         0.2    setosa
...
```

⁴ \mathbf{S} será definida positiva si $n \geq p + 1$.

```
# extraemos solamente variables numéricas
> x <- iris[,1:4]
> pairs(x, col = iris$Species) # Fig.3
```

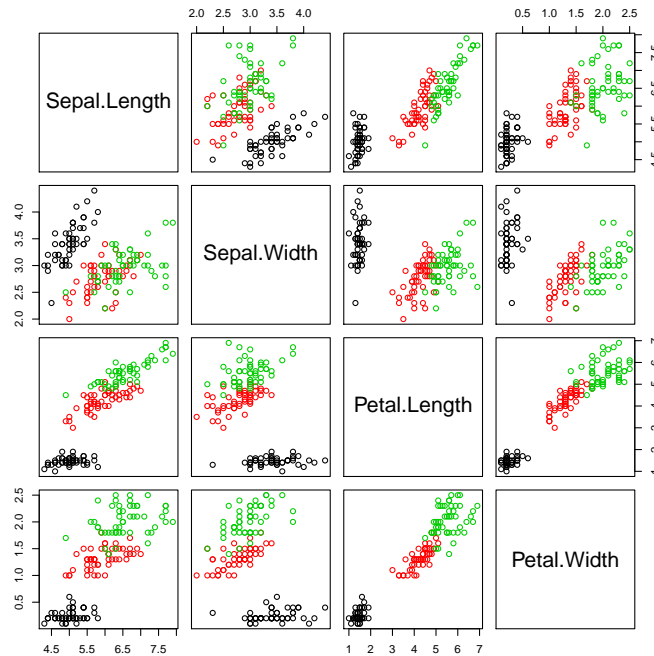


Figura 3. Gráficos de dispersión para los datos de Iris, según variedades Setosa (negro), Versicolor (rojo) y Virginica (verde).

El gráfico de dispersión desplegado en la Figura 3 revela una separación evidente en dos grupos (Setosa versus Versicolor y Virginica). Los siguientes comandos invocan funciones desde la biblioteca `fastmatrix`

```
# Carga biblioteca 'fastmatrix'
> library(fastmatrix)
# Cálculo de estadísticas de resumen multivariadas
> z <- cov.weighted(x) # por defecto, los 'pesos' son 1
> xbar <- z$mean
> S <- z$cov
> R <- cov2cor(z$cov)
> b1 <- skewness(x)
> b2 <- kurtosis(x)

# análogamente podemos calcular S y R usando:
> S <- cov(x)
> R <- cor(x)
```

En la línea de comandos de R podemos escribir `xbar`, `R`, `b1` y `b2` para obtener, el vector de medias, la matriz de correlación y el coeficiente de sesgo y curtosis muestrales, dados por:

$$\bar{x} = \begin{pmatrix} 5.8433 \\ 3.0573 \\ 3.7580 \\ 1.1993 \end{pmatrix}, \quad R = \begin{pmatrix} 1.0000 & -0.1176 & 0.8718 & 0.8179 \\ -0.1176 & 1.0000 & -0.4284 & -0.3661 \\ 0.8718 & -0.4284 & 1.0000 & 0.9629 \\ 0.8179 & -0.3661 & 0.9629 & 1.0000 \end{pmatrix},$$

mientras que $b_{1p} = 2.6972$ y $b_{2p} = 23.7397$, respectivamente.

1.5. Regresión lineal simple

Ahora estamos enfocados en situaciones tales que la variable x permite predecir o explicar la respuesta y . En regresión lineal tenemos el modelo,

$$y_i = \alpha + \beta x_i + \epsilon_i, \quad i = 1, \dots, n. \quad (1.28)$$

El objetivo es, basado en el conjunto de datos pareados, $(x_1, y_1), \dots, (x_n, y_n)$ determinar α y β tal que produzcan el mejor ajuste. Utilizaremos el *método de mínimos cuadrados ordinarios* (OLS) definido como:

$$\min_{\theta} S(\theta),$$

con $\theta = (\alpha, \beta)^\top$ y

$$S(\theta) = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 = \sum_{i=1}^n \epsilon_i^2.$$

OBSERVACIÓN 1.61. OLS minimiza las distancias *verticales* a la recta de regresión. Es decir, esto refleja que los x_i 's están *fijados* (o bien que se asumen conocidos).⁵

La función $S(\theta)$ también es conocida como suma de cuadrados de los errores. De este modo, debemos determinar $\theta = (\alpha, \beta)^\top$, mediante resolver las condiciones de primer orden,

$$\begin{aligned} \frac{\partial S(\theta)}{\partial \alpha} &= -2 \sum_{i=1}^n (y_i - \alpha - \beta x_i) = 0, \\ \frac{\partial S(\theta)}{\partial \beta} &= -2 \sum_{i=1}^n x_i (y_i - \alpha - \beta x_i) = 0, \end{aligned}$$

esto lleva a las *ecuaciones normales*:

$$\begin{aligned} \sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta} x_i) &= \sum_{i=1}^n y_i - n\hat{\alpha} - \hat{\beta} \sum_{i=1}^n x_i = 0, \\ \sum_{i=1}^n x_i (y_i - \hat{\alpha} - \hat{\beta} x_i) &= \sum_{i=1}^n x_i y_i - \hat{\alpha} \sum_{i=1}^n x_i - \hat{\beta} \sum_{i=1}^n x_i^2 = 0. \end{aligned}$$

Multiplicando la primera ecuación por n^{-1} y resolviendo para $\hat{\alpha}$, obtenemos

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}.$$

⁵Una alternativa es usar mínimos cuadrados totales o modelos con errores en las variables.

Substituyendo $\hat{\alpha}$ en la segunda ecuación, tenemos

$$\begin{aligned} \sum_{i=1}^n x_i y_i - (\bar{y} - \hat{\beta} \bar{x}) \sum_{i=1}^n x_i - \hat{\beta} \sum_{i=1}^n x_i^2 &= 0 \\ \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} + \hat{\beta} n \bar{x}^2 - \hat{\beta} \sum_{i=1}^n x_i^2 &= 0 \\ \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} - \hat{\beta} \left(\sum_{i=1}^n x_i^2 - n \bar{x}^2 \right) &= 0 \\ \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) - \hat{\beta} \sum_{i=1}^n (x_i - \bar{x})^2 &= 0. \end{aligned}$$

Lo que lleva a

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{\text{var}(\mathbf{x})}.$$

De este modo, la *recta de regresión* asume la forma:

$$\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i, \quad i = 1, \dots, n,$$

y llamamos a \hat{y}_i el valor “*predicho*” o “*ajuste*” para x_i . Se define

$$e_i = y_i - \hat{y}_i = y_i - \hat{\alpha} - \hat{\beta} x_i, \quad i = 1, \dots, n,$$

como el i -ésimo *residuo*, mientras que una medida de variabilidad es dada por:

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{\alpha} - \hat{\beta} x_i)^2 = \frac{1}{n-2} S(\hat{\boldsymbol{\theta}}),$$

con $\hat{\boldsymbol{\theta}} = (\hat{\alpha}, \hat{\beta})^\top$. Debemos resaltar que, el promedio de los valores predichos $\hat{y}_1, \dots, \hat{y}_n$, asume la forma

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \hat{y}_i &= \frac{1}{n} \sum_{i=1}^n (\hat{\alpha} + \hat{\beta} x_i) = \hat{\alpha} + \hat{\beta} \bar{x} \\ &= (\bar{y} - \hat{\beta} \bar{x}) + \hat{\beta} \bar{x} = \bar{y}. \end{aligned}$$

Haciendo $R = \text{cor}(\mathbf{x}, \mathbf{y})$, sigue que

$$\hat{y}_i - \bar{y} = \hat{\alpha} + \hat{\beta} x_i - \bar{y} = \hat{\beta} (x_i - \bar{x}) = R \frac{s_y}{s_x} (x_i - \bar{x}),$$

es decir,

$$(\hat{y}_i - \bar{y})^2 = R^2 \frac{s_y^2}{s_x^2} (x_i - \bar{x})^2.$$

Lo que nos lleva a escribir,

$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = R^2 \frac{s_y^2}{s_x^2} \sum_{i=1}^n (x_i - \bar{x})^2 = R^2 \frac{s_y^2}{s_x^2} s_x^2 = R^2 s_y^2 = R^2 \sum_{i=1}^n (y_i - \bar{y})^2.$$

A partir de esta ecuación, sigue que

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{s_{\text{AJUSTE}}^2}{s_{\text{DATOS}}^2},$$

que se denomina *coeficiente de determinación*. Usando además que

$$y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y}),$$

podemos mostrar que

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2,$$

es decir,

$$s_{\text{DATOS}}^2 = s_{\text{RESIDUOS}}^2 + s_{\text{AJUSTE}}^2.$$

Como $s_{\text{RESIDUOS}}^2 = S(\hat{\theta})$, tenemos que $s_{\text{AJUSTE}}^2 = R^2 s_{\text{DATOS}}^2$. Luego, sigue que:

$$s_{\text{RESIDUOS}}^2 = (1 - R^2) s_{\text{DATOS}}^2.$$

Esto permite interpretar R^2 como la proporción de varianza de los datos que puede ser explicada por el ajuste. En efecto, $0 \leq R^2 \leq 1$ permite medir la calidad (o bondad) del ajuste.

EJEMPLO 1.62. Debemos ser precavidos al llevar un análisis de regresión y confiar de medidas tales como el R^2 para evaluar la calidad del ajuste. Para notar este tipo de situaciones, considere los datos introducidos por [Anscombe \(1973\)](#). Este corresponde a cuatro conjunto de datos con 11 observaciones que aunque son bastante diferentes tienen estadísticas de resumen idénticas (promedio, varianza, correlación, coeficientes de regresión, etc.). En efecto, considere la siguiente figura

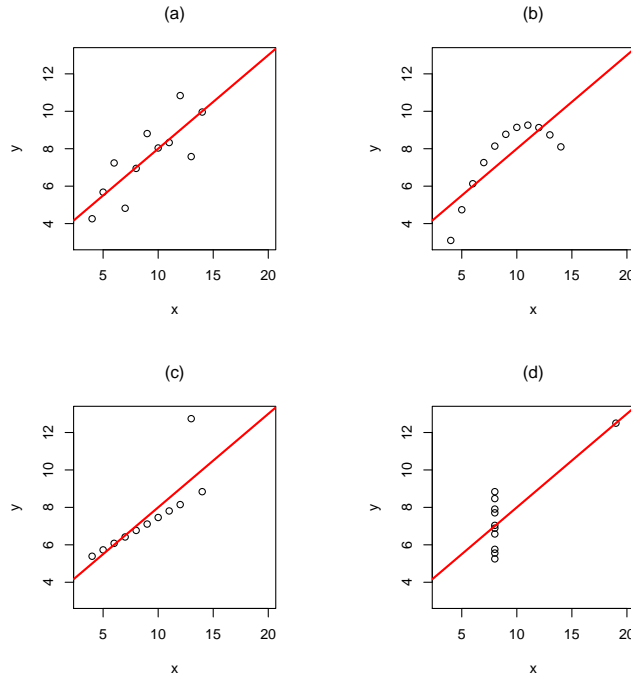


Figura 4. Cuarteto de regresiones idénticas de [Anscombe \(1973\)](#).

Por ejemplo, en la Figura 4(b), los datos presentan una tendencia cuadrática, mientras que en los Paneles (c) y (d) se aprecia el efecto de observaciones atípicas sobre la recta de regresión. Más aún, eliminando la observación extrema en el conjunto de datos en la Figura 4(d), la regresión deja de tener sentido. Para obtener los resultados del ajuste mediante regresión por mínimos cuadrados podemos considerar los siguientes comandos de R

```
# datos de Anscombe
> anscombe
  x1 x2 x3 x4    y1    y2    y3    y4
1  10 10 10  8  8.04 9.14  7.46  6.58
2   8  8  8  8  6.95 8.14  6.77  5.76
3  13 13 13  8  7.58 8.74 12.74  7.71
...
10  7  7  7  8  4.82 7.26  6.42  7.91
11  5  5  5  8  5.68 4.74  5.73  6.89

# Ajuste mediante OLS (para cada uno de los modelos)
> fm1 <- lm(y1 ~ x1, data = anscombe)
> fm2 <- lm(y2 ~ x2, data = anscombe)
> fm3 <- lm(y3 ~ x3, data = anscombe)
> fm4 <- lm(y4 ~ x4, data = anscombe)

# Figura 4.a (otros paneles son análogos)
> par(pty = "s")
> plot(y1 ~ x1, data = anscombe, xlim = c(4,20),
+      ylim = c(3,13), xlab = "x", ylab = "y")
> abline(coef(fm1), col = "red", lwd = 2) # línea en rojo

# Salida de resultados:
> summary(fm1)
Call:
lm(formula = y1 ~ x1, data = anscombe)

Residuals:
    Min       1Q   Median       3Q      Max
-1.92127 -0.45577 -0.04136  0.70941  1.83882

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   3.0001     1.1247   2.667  0.02573 *
x1             0.5001     0.1179   4.241  0.00217 **

Residual standard error: 1.237 on 9 degrees of freedom
Multiple R-squared:  0.6665, Adjusted R-squared:  0.6295
F-statistic: 17.99 on 1 and 9 DF, p-value: 0.00217

# Almacena coeficientes, residuos y valores predichos
> cf <- coef(fm1)
> res <- resid(fm1)
> fit <- fitted(lm1)
```

```
# Figura 5.a (otros paneles son análogos)
> par(pty = "s")
> plot(fit, res, xlab="Valores predichos", ylab="Residuos")
> abline(h = 0, lty = 2, col = "gray", lwd = 2) # línea gris
```

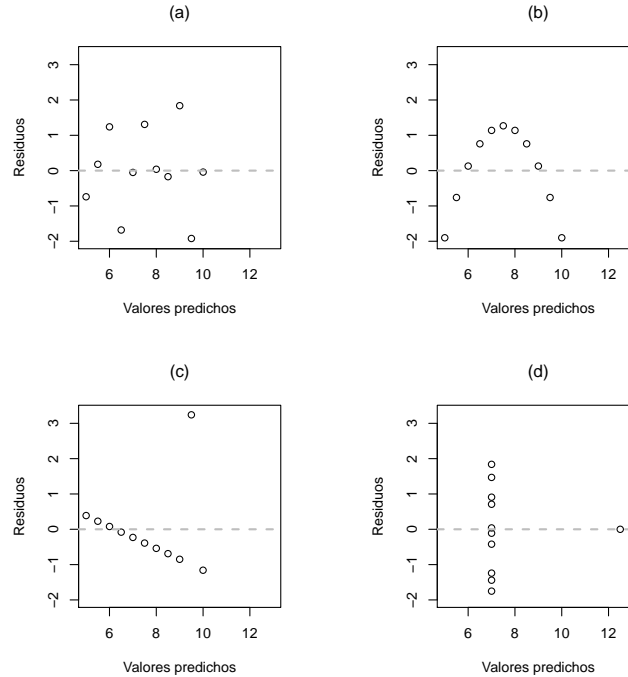


Figura 5. Gráfico de residuos para el cuarteto de regresiones idénticas de [Anscombe \(1973\)](#).

Para el modelo en Ecuación (1.28), es posible mostrar que:

$$\sum_{i=1}^n e_i = 0, \quad \text{y} \quad \sum_{i=1}^n e_i \hat{y}_i = 0.$$

De este modo, el gráfico de dispersión de residuos contra valores predichos dado en la Figura 5 no debería presentar algún tipo de comportamiento sistemático. Es decir, el único conjunto de datos donde el modelo de regresión lineal parece razonable es el desplegado en el Panel (a) de la Figura 4. En efecto, el gráfico de residuos en la Figura 5 (a) es el único que no presenta algún tipo de tendencia. En cambio, se debe incluir un término cuadrático al modelo de la Figura 4 (b). Mientras que el rol de las observaciones atípicas visibles en los Paneles (c) y (d) de la Figura 4 se aprecia claramente en los gráficos de residuos. Es decir, corresponden a un outlier (observación alejada en la ‘respuesta’) y un leverage (observación alejada en los ‘regresores’), respectivamente.

Debemos destacar que para este ejemplo sencillo es bastante fácil determinar que el modelo lineal no es apropiado para las configuraciones de datos en los Paneles

(b), (c) y (d) presentados en la Figura 4. En general, no resulta fácil identificar situaciones donde se aprecie que el modelo propuesto está mal ajustado. Este tipo de situaciones es más difícil de determinar conforme se tiene de un mayor volumen de datos, o bien la cantidad de regresores aumenta.⁶

1.6. Resúmenes gráficos

Para introducir ideas consideremos una imagen digital de tamaño $r \times c$. Una imagen digital es un arreglo de tamaño $r \times c$ que puede ser conceptualizada como una matriz de datos $\mathbf{X} = (X_{ij})$, para $i = 1, \dots, r$, $j = 1, \dots, c$, donde el valor de la variable X en la posición (i, j) corresponde a un valor en una escala de grises en el intervalo $[0, 1]$ tal que 0 representa el color negro y 1 representa el blanco. De este modo, $\{X_{ij} : 1 \leq i \leq r, 1 \leq j \leq c\}$ representa el conjunto de intensidades de gris en cada una de las $r \times c$ posiciones en el espacio bidimensional. Cada posición (i, j) es denominada un *pixel* (picture element). La Figura 6 presenta una imagen muy popular en ingeniería, específicamente en el área de tratamiento de imágenes, llamada frecuentemente *Lenna* y fue popularizada hace varias décadas atrás.⁷



Figura 6. Imagen *Lenna* de tamaño 512×512 .

Podemos considerar una imagen como un vector $\mathbf{x} = \text{vec}(\mathbf{X})$ en el que hemos concatenado las columnas de \mathbf{X} . Para la imagen Lenna tenemos $n = 512 \times 512 = 262144$ observaciones. Un resumen para el conjunto de datos *Lenna* es obtenido usando los siguientes comandos en R:

```
# carga datos de 'Lenna' desde URL
> lena <- "http://fosorios.mat.utfsm.cl/files/data/lenna.png"
> download.file(lena, z, mode = "wb")
> library(png) # biblioteca para cargar imágenes PNG
> lena <- readPNG(z)[, , 1] # sólo el 1er canal es necesario
```

⁶En la asignatura *Análisis de Regresión* se revisará el aspecto de la crítica del modelo con mayor profundidad.

⁷Para más detalles sobre esta imagen, consultar la página web: <http://www.lenna.org/>

```

# carga biblioteca 'SpatialPack'
> library(SpatialPack)
> lena <- normalize(lena)
> plot(as.raster(lena)) # Figura 6

# convierte en un vector de datos
> x <- as.vector(lena)
> summary(x)
      Min. 1st Qu.  Median    Mean 3rd Qu.   Max.
0.00000  0.2712  0.4443  0.4316  0.5809  1.0000

# cuantiles más utilizados
> qx <- quantile(x)
> qx
      0%      25%      50%      75%     100%
0.0000000 0.2712455 0.4442928 0.5808888 1.0000000

> qx <- as.vector(qx)[-c(1,5)] # remueve extremos
> IQR <- qx[3] - qx[1]
> IQR # rango intercuartílico
[1] 0.3096433

```

Podemos notar que, debido a la normalización de los datos, el valor mínimo y máximo corresponden a 0 y 1, respectivamente, mientras que el rango intercuartílico, esto es, la distancia en la que se concentra el 50% central de la muestra resulta $IQR = Q_3 - Q_1 = 0.5808 - 0.2712 = 0.3096$.

A continuación se presenta dos herramientas gráficas para una muestra x_1, x_2, \dots, x_n que permiten visualizar de forma simple la distribución de una variable de interés X , conocidas como *histograma* y *boxplot* (o diagrama de cajón-con-bigotes).

1.6.1. Histograma. Un histograma es una aproximación a la densidad desconocida construido a partir de los datos observados. Considere la siguiente definición.

DEFINICIÓN 1.63. Sea x_1, x_2, \dots, x_n una muestra de n observaciones y suponga que los datos son subdivididos en k intervalos, digamos I_1, I_2, \dots, I_k . Entonces, la frecuencia absoluta de la clase I_j es la cantidad de observaciones n_j de la muestra, que pertenecen al intervalo I_j . Evidentemente tenemos que:

$$n_j \geq 0, \quad \text{y} \quad \sum_{j=1}^k n_j = n.$$

La frecuencia relativa asociada a la clase I_j se define como:

$$f_j = \frac{n_j}{n}, \quad j = 1, \dots, k.$$

En este caso es fácil ver que $f_j \geq 0$ y $\sum_{j=1}^k f_j = 1$. De este modo, un *histograma* es un gráfico de f_j (o n_j) versus I_j .

Un histograma es un diagrama de frecuencias y resume la cantidad de observaciones por unidad de longitud. Luego, este diagrama permite visualizar la distribución de la variable de interés, tal como se presenta en la Figura 7. En este caso, ambos

histogramas tienen 14 clases. Este corresponde a un parámetro gráfico que puede modificarse para adaptarse a la estructura de la muestra que se tiene disponible. En algunos casos se grafica f_j/h versus I_j , donde h es la amplitud de los intervalos I_j . Este gráfico tiene la particularidad de que el área bajo la curva es dada por

$$A = \sum_{j=1}^k A_j = \sum_{j=1}^k (f_j/h) \cdot h = \sum_{j=1}^k f_j = 1.$$

En efecto, resultará claro más adelante que esta propiedad está asociada al concepto de probabilidad.

Debemos enfatizar que un histograma es una herramienta descriptiva. Más adelante estudiaremos algunas curvas llamadas *funciones de densidad de probabilidad* asociada a ciertas variables de interés. Sin embargo, histogramas han sido criticados por tener algunas inconvenientes notables. Además de la necesidad de escoger el número de intervalos k (o equivalentemente la amplitud de los intervalos h), frecuentemente pueden resultar engañosos si se interpretan en exceso.

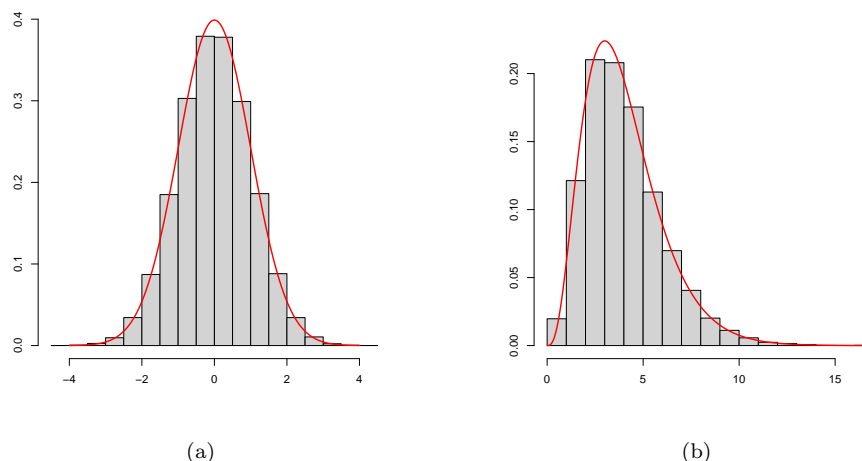


Figura 7. Histograma y densidad teórica para una muestra de 10 000 observaciones provenientes desde una distribución (a) normal y (b) gama.

1.6.2. Boxplot (Diagrama de cajón con bigotes). Este tipo de diagrama corresponde a un tipo diferente de visualización que nos permite explorar la ubicación, la escala, la asimetría y las colas de una densidad así como la presencia de observaciones atípicas (outliers). En contraste con el histograma el boxplot ofrece una descripción mucho más gruesa de la estructura de la muestra y no requiere la especificación de parámetros para la calibración del gráfico. Esta representación usualmente se realiza en forma de una caja, lo que explica el nombre del gráfico. A continuación se presenta el procedimiento para su construcción.

DEFINICIÓN 1.64. Considere x_1, x_2, \dots, x_n una muestra de n observaciones y sea:

- Q_1 , me y Q_3 respectivamente, el primer cuartil, la mediana y el tercer cuartil de $\{x_1, \dots, x_n\}$.
- Cálculo de los “bigotes” W_1 y W_2 , como:

$$W_1 = \min_{1 \leq j \leq n} \{x_j : x_j \geq Q_1 - 1.5IQR\}$$

$$W_2 = \max_{1 \leq j \leq n} \{x_j : x_j \leq Q_3 + 1.5IQR\}$$

- $O = \{i \in \{1, 2, \dots, n\} : x_i \notin [W_1, W_2]\}$.

De este modo, el *boxplot* de x_1, \dots, x_n es el registro de W_1, Q_1, me, Q_3, W_2 y $\{x_j : j \in O\}$ sobre la recta real.

Aunque la definición puede ser un tanto difícil de visualizar, los gráficos en la Figura 8 permiten comprender la idea. La línea central corresponde a la mediana, mientras que los extremos del cajón representan Q_1 y Q_3 . Note que W_1 y W_2 representan los extremos del bigote y las observaciones atípicas son indicados por puntos. En el Panel (a) tenemos una distribución simétrica en torno de cero, mientras que la distribución en el Panel (b) presenta asimetría positiva.

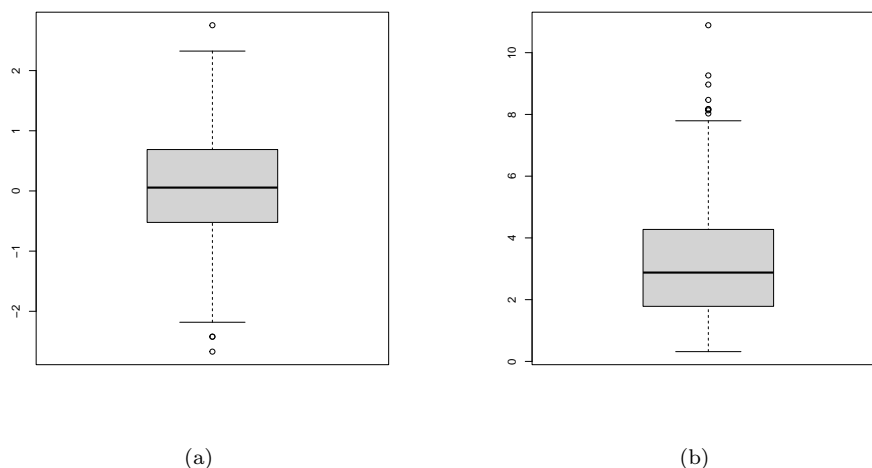


Figura 8. Boxplot para una muestra de 300 observaciones provenientes desde una distribución (a) normal y (b) gama.

Un aspecto destacable del boxplot, es que permite la comparación de varias muestras permitiendo la visualización simultánea de medidas de tendencia central y dispersión mediante construir una secuencia de gráficos de cajón-con-bigotes.

Retomando el conjunto de datos de la imagen Lenna, el siguiente código en R, permite obtener el histograma, la curva de densidad estimada y el boxplot:

```
# histograma y densidad estimada (Fig. 9)
> hist(x, freq = FALSE, main = "", xlab = "Lenna",
+      ylab = "Densidad")
```

```

> plot(density(x), main = "", xlab = "bandwith = 0.0161",
+      ylab = "Densidad")
> boxplot(x)
> xbar <- mean(x) # promedio
> abline(h = xbar, lty = 2, col = "red") # línea segmentada

```

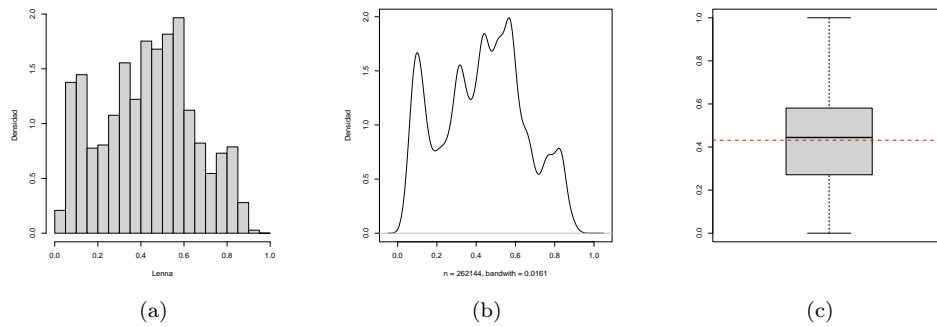


Figura 9. Histograma, densidad estimada y boxplot para los datos de la imagen Lenna.

Es interesante destacar que imágenes en color se suelen representar como un arreglo con tres canales, usualmente rojo (R), verde (G) y azul (B), lo que son combinados de manera apropiada para producir una representación a color. Considere la imagen Lenna en formato a color,



Figura 10. Imagen Lenna en versión a color.

De este modo, podemos considerar que las intensidades de rojo, verde y azul, representan 3 conjuntos de datos. Con fines comparativos, Figuras 11 y 12 presentan histogramas y boxplots para cada uno de los canales RGB de la imagen Lenna.

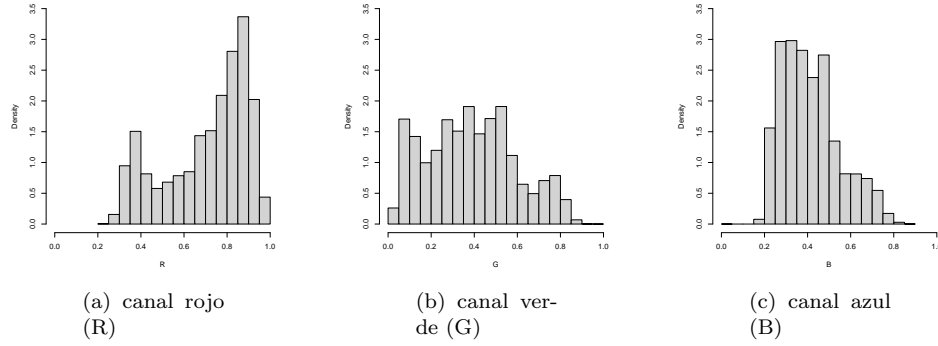


Figura 11. Histogramas para cada uno de los canales RGB de los datos de la imagen Lenna.

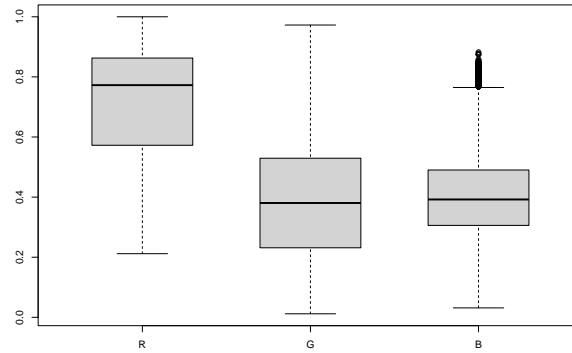


Figura 12. Boxplots para cada uno de los canales RGB de los datos de la imagen Lenna.

Es notable como los boxplots en la Figura 12 nos permite hacer una comparación rápida de la configuración de cada uno de los canales RGB.

El campo de visualización es tópico amplio con desafíos muy relevantes dado la demanda por resumir grandes volúmenes de información usando un despliegue gráfico. Para mayores detalles sobre visualización así como por técnicas de análisis exploratorio consultar [Tukey \(1977\)](#), [Cleveland \(1993\)](#) y [Wilkinson \(1999\)](#).

Nociones de Probabilidad

2.1. Preliminares

DEFINICIÓN 2.1. El conjunto Ω , de todos los resultados posibles de un experimento aleatorio es llamado *espacio muestral*.

EJEMPLO 2.2. Suponga que lanzamos de una moneda. De este modo, tenemos sólo dos resultados posibles

$$\Omega = \{C, S\}.$$

EJEMPLO 2.3. Para el caso de las notas de la asignatura MAT-021 para un grupo de estudiantes escogidos al azar, sigue que

$$\Omega = \{0, 1, \dots, 99, 100\}.$$

EJEMPLO 2.4. El espacio muestral para el tiempo de duración de un artículo es dado por:

$$\Omega = [0, \infty).$$

DEFINICIÓN 2.5. Un evento (o suceso) es cualquier colección de resultados posibles de un experimento aleatorio, esto es cualquier subconjunto de Ω (incluyendo al propio Ω).

OBSERVACIÓN 2.6. Sea $A \subset \Omega$, diremos que A ocurre si $\omega \in A$ con $\omega \in \Omega$ es un resultado asociado a un experimento aleatorio. Evidentemente, $\omega \notin A$ si y sólo si A no ocurre.

A continuación describimos algunas operaciones sobre conjuntos que usaremos frecuentemente:

$$A \subset B \iff x \in A \implies x \in B.$$

$$A = B \iff A \subset B \text{ y } B \subset A.$$

Unión: La unión entre A y B denotada por $A \cup B$ es definida como:

$$A \cup B = \{x : x \in A, \text{ o } x \in B\}.$$

Intersección: La intersección entre A y B escrita como $A \cap B$ se define como:

$$A \cap B = \{x : x \in A, \text{ y } x \in B\}.$$

Complemento: El complemento de A , denotado como A^c es el conjunto de todos los elementos que *no* están en A :

$$A^c = \{x : x \notin A\}.$$

PROPIEDAD 2.7. Sean A, B y C tres eventos definidos en Ω . Tenemos:

$$(a) \quad A \cup B = B \cup A, \quad A \cap B = B \cap A.$$

- (b) $A \cup (B \cap C) = (A \cup B) \cap C$,
 $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$.
- (c) $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$,
 $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$.
- (d) $(A \cup B)^c = A^c \cap B^c$, y $(A \cap B)^c = A^c \cup B^c$.

La Propiedad 2.7 (d) es conocida frecuentemente como *Leyes de De Morgan*. La demostración de cada una de estas propiedades se dejan como ejercicio para el lector.

Note que las operaciones de unión e intersección se pueden extender a colecciones infinitas de conjuntos,

$$\bigcup_{i=1}^{\infty} A_i = \{x \in \Omega : x \in A_i, \text{ para algún } i\}.$$

$$\bigcap_{i=1}^{\infty} A_i = \{x \in \Omega : x \in A_i, \text{ para todo } i\}.$$

DEFINICIÓN 2.8. Dos eventos A y B son disjuntos (o excluyentes) si

$$A \cap B = \emptyset.$$

Los eventos A_1, A_2, \dots , son disjuntos por pares (o mutuamente excluyentes) si

$$A_i \cap A_j = \emptyset, \quad i \neq j.$$

DEFINICIÓN 2.9. Si A_1, A_2, \dots son disjuntos por pares y $\bigcup_{i=1}^{\infty} A_i = \Omega$, entonces la colección se denomina una *partición* de Ω .

EJEMPLO 2.10. Los conjuntos $A_i = [i, i+1)$, para $i = 0, 1, 2, \dots$ forman una partición de $[0, \infty)$.

EJEMPLO 2.11. $\Omega = A \cup A^c$ es una partición.

DEFINICIÓN 2.12. Sea A, B dos sucesos, diremos que A implica B si y sólo si $A \subseteq B$.

2.2. Conceptos básicos

Para todo evento $A \subset \Omega$ deseamos asociar un número entre cero y uno llamado probabilidad de A .

DEFINICIÓN 2.13. Una colección de subconjuntos de Ω es llamado σ -álgebra y es denotada por \mathcal{F} si satisface las propiedades:

- (a) $\emptyset \in \mathcal{F}$.
- (b) Si $A \in \mathcal{F} \implies A^c \in \mathcal{F}$.
- (c) Si $A_1, A_2, \dots \in \mathcal{F} \implies \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$

Note que $\emptyset \subset \Omega$ y $\Omega = \emptyset^c$, así por la Propiedad (a) y (b) sigue que $\Omega \in \mathcal{F}$.

Además, si $A_1, A_2, \dots \in \mathcal{F}$ entonces $A_1^c, A_2^c, \dots \in \mathcal{F}$ y de este modo, $\bigcup_{i=1}^{\infty} A_i^c \in \mathcal{F}$. Por las leyes de De Morgan, tenemos

$$\left(\bigcup_{i=1}^{\infty} A_i^c \right)^c = \bigcap_{i=1}^{\infty} A_i.$$

Es decir, tenemos que $\bigcap_{i=1}^{\infty} A_i \in \mathcal{F}$.

OBSERVACIÓN 2.14. Asociado a un espacio muestral Ω puede haber muchas σ -álgebras. Por ejemplo, la colección $\{\emptyset, \Omega\}$ es σ -álgebra (minimal).

DEFINICIÓN 2.15. Sea Ω un espacio muestral y sea \mathcal{F} un σ -álgebra, decimos que el par (Ω, \mathcal{F}) es un espacio medible. Si $A \in \mathcal{F}$ decimos que A es medible.

DEFINICIÓN 2.16 (Probabilidad). Dado un espacio muestral Ω y un σ -álgebra asociada \mathcal{F} , una función de probabilidad P , $P : \mathcal{F} \rightarrow \mathbb{R}$, satisface:

- (a) $P(A) \geq 0$, para todo $A \in \mathcal{F}$.
- (b) $P(\Omega) = 1$.
- (c) Si A_1, A_2, \dots son disjuntos por pares (mutuamente excluyentes), entonces

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i).$$

Sea (Ω, \mathcal{F}) un espacio de medida y P una medida de probabilidad definida en \mathcal{F} . Entonces (Ω, \mathcal{F}, P) se denomina *espacio de probabilidad*.

RESULTADO 2.17. Si P es una función de probabilidad y A es cualquier conjunto en \mathcal{F} , entonces

- (a) $P(\emptyset) = 0$.
- (b) $P(A) \leq 1$.
- (c) $P(A^c) = 1 - P(A)$.

DEMOSTRACIÓN. Primero considere (c). Como A y A^c son una partición de Ω , sigue que

$$P(A \cup A^c) = P(\Omega) = 1,$$

además $A \cap A^c = \emptyset$ son disjuntos, luego

$$P(A \cup A^c) = P(A) + P(A^c) = 1,$$

lo que permite mostrar el resultado. Como $P(A^c) \geq 0$, (b) sigue desde (c). Finalmente, para probar (a) note que $\Omega = \Omega \cup \emptyset$ y como $\Omega \cap \emptyset = \emptyset$, tenemos

$$1 = P(\Omega) = P(\Omega \cup \emptyset) = P(\Omega) + P(\emptyset),$$

de este modo $P(\emptyset) = 0$. □

RESULTADO 2.18. Si P es una función de probabilidad y $A, B \in \mathcal{F}$. Entonces,

- (a) $P(B \cap A^c) = P(B) - P(A \cap B)$.
- (b) $P(A \cup B) = P(A) + P(B) - P(A \cap B)$.
- (c) Si $A \subseteq B \implies P(A) \leq P(B)$.

DEMOSTRACIÓN. Para notar (a) considere que para A y B conjuntos cualquiera

$$B = (B \cap A) \cup (B \cap A^c)$$

(pues $B = B \cup (A \cap A^c) = B \cup \emptyset$). Luego,

$$P(B) = P(\{B \cap A\} \cup \{B \cap A^c\}) = P(B \cap A) + P(B \cap A^c).$$

En efecto, $(B \cap A) \cap (B \cap A^c) = B \cap (A \cap A^c) \cap B = \emptyset$. Para probar (b), note que

$$(A \cup B) = A \cup (B \cap A^c).$$

Además,

$$A \cap (B \cap A^c) = (A \cap A^c) \cap B = \emptyset \cap B = \emptyset.$$

Por tanto, sigue que

$$P(A \cup B) = P(A) + P(B \cap A^c) = P(A) + P(B) - P(A \cap B),$$

por parte (a).

Si $A \subseteq B$ entonces $A \cap B = A$. De este modo, usando (a) tenemos

$$0 \leq P(B \cap A^c) = P(B) - P(A),$$

y (c) es verificado. □

Como $P(A \cup B) \leq 1$, reagrupando términos tenemos

$$P(A \cap B) \geq P(A) + P(B) - 1,$$

la desigualdad anterior es un caso particular de la desigualdad de Bonferroni.

RESULTADO 2.19. Si P es una función de probabilidad. Entonces,

- (a) $P(A) = \sum_{i=1}^{\infty} P(A \cap C_i)$ para cualquier partición C_1, C_2, \dots .
- (b) $P(\bigcup_{i=1}^{\infty} A_i) \leq \sum_{i=1}^{\infty} P(A_i)$ para conjuntos A_1, A_2, \dots cualquiera.¹

DEMOSTRACIÓN. Dado que C_1, C_2, \dots forman una partición, tenemos $C_i \cap C_j = \emptyset$ para todo $i \neq j$ y $\Omega = \bigcup_{i=1}^{\infty} C_i$. De ahí que

$$A = A \cap \Omega = A \cap \left(\bigcup_{i=1}^{\infty} C_i \right) = \bigcup_{i=1}^{\infty} (A \cap C_i).$$

De este modo,

$$P(A) = P\left(\bigcup_{i=1}^{\infty} (A \cap C_i)\right) = \sum_{i=1}^{\infty} P(A \cap C_i),$$

pues, dado que los C_i 's son disjuntos, también lo es la secuencia $\{A \cap C_i\}_{i=1}^{\infty}$.

Para establecer (b) se construye una colección disjunta A_1^*, A_2^*, \dots tal que

$$\bigcup_{i=1}^{\infty} A_i^* = \bigcup_{i=1}^{\infty} A_i.$$

Defina A_i^* , como

$$\begin{aligned} A_1^* &= A_1, \\ A_i^* &= A_i - \left(\bigcup_{j=1}^{i-1} A_j \right), \quad i = 2, 3, \dots, \end{aligned}$$

donde $A - B = A \cap B^c$ denota la diferencia entre conjuntos. Podemos notar que $\bigcup_{i=1}^{\infty} A_i^* = \bigcup_{i=1}^{\infty} A_i$, luego

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = P\left(\bigcup_{i=1}^{\infty} A_i^*\right) = \sum_{i=1}^{\infty} P(A_i^*),$$

¹Esta es conocida como la *Desigualdad de Boole*.

pues los A_i^* 's son disjuntos. En efecto,

$$\begin{aligned} A_i^* \cap A_k^* &= \left\{ A_i - \left(\bigcup_{j=1}^{i-1} A_j \right) \right\} \cap \left\{ A_k - \left(\bigcup_{j=1}^{k-1} A_j \right) \right\} \\ &= \left\{ A_i \cap \left(\bigcup_{j=1}^{i-1} A_j \right)^c \right\} \cap \left\{ A_k \cap \left(\bigcup_{j=1}^{k-1} A_j \right)^c \right\} \\ &= \left\{ A_i \cap \left(\bigcap_{j=1}^{i-1} A_j^c \right) \right\} \cap \left\{ A_k \cap \left(\bigcap_{j=1}^{k-1} A_j^c \right) \right\}, \end{aligned}$$

si $i > k$ entonces la primera intersección está contenida en A_k^c , luego esa intersección será vacía. Si $i < k$ el argumento es similar. Además, por construcción $A_i^* \subset A_i$ de modo que $P(A_i^*) \leq P(A_i)$. Por tanto,

$$\sum_{i=1}^{\infty} P(A_i^*) \leq \sum_{i=1}^{\infty} P(A_i),$$

estableciendo el resultado. \square

OBSERVACIÓN 2.20. Usando la desigualdad de Boole, tenemos

$$P\left(\bigcup_{i=1}^n A_i^c\right) \leq \sum_{i=1}^n P(A_i^c),$$

y como $\bigcup A_i^c = (\bigcap A_i)^c$ y $P(A_i^c) = 1 - P(A_i)$, tenemos

$$1 - P\left(\bigcap_{i=1}^n A_i\right) \leq n - \sum_{i=1}^n P(A_i),$$

es decir

$$P\left(\bigcap_{i=1}^n A_i\right) \geq \sum_{i=1}^n P(A_i) - (n - 1),$$

que es conocida como la desigualdad de Bonferroni.

2.3. Espacios muestrales finitos

En esta sección consideraremos que

$$\Omega = \{\omega_1, \omega_2, \dots, \omega_k\},$$

para caracterizar $P(A)$ supondremos eventos elementales, es decir $A = \{\omega_i\}$ y definimos $p_i = P(\{\omega_i\})$ la probabilidad de $\{\omega_i\}$ tal que,

- (a) $p_i \geq 0$, $i = 1, \dots, k$.
- (b) $p_1 + p_2 + \dots + p_k = 1$.

Suponga ahora que $A = \{\omega_{j_1}, \dots, \omega_{j_r}\}$ está formado por r elementos de Ω , luego

$$P(A) = p_{j_1} + \dots + p_{j_r}.$$

Adicionalmente, supondremos que cada $\{\omega_i\}$ es igualmente probable. Entonces,

$$p_i = P(\{\omega_i\}) = \frac{1}{k}.$$

Luego, para un evento $A = \{\omega_{j_1}, \dots, \omega_{j_r}\}$ sigue que

$$P(A) = \frac{r}{k},$$

o bien,

$$P(A) = \frac{\#(A)}{\#(\Omega)},$$

donde $\#(A)$ denota la cardinalidad del conjunto A . Debemos resaltar que esta *no* es una definición general, sino que apropiada *sólo* bajo el supuesto de espacios muestrales finitos y equiprobables.

EJEMPLO 2.21. Considere un lote de 100 artículos con 20 defectuosos y 80 no defectuosos. Se eligen 10 artículos al azar, sin substituir un artículo antes de elegir el próximo. ¿Cuál es la probabilidad de que exactamente la mitad de los artículos escogidos sean defectuosos? Note que el espacio muestral consta de todos los vectores $(i_1, i_2, \dots, i_{10})$ extraídos desde el lote de 100 artículos ¿Cómo contar cuántos son? ¿Cuáles satisfacen la condición que define $P(A)$?

2.4. Técnicas de conteo

El objetivo es determinar el número de subconjuntos distintos que se pueden formar con un conjunto de elementos dados.

2.4.1. Principio multiplicativo. Suponga que se desarrolla un experimento y digamos que podemos subdividir este experimento en 2 etapas. Si la primera etapa puede desarrollarse de n_1 maneras y la segunda etapa en n_2 manera, entonces el experimento puede llevarse a cabo en

$$n_1 \cdot n_2,$$

maneras.

OBSERVACIÓN 2.22. Esto puede extenderse a un experimento que se desarrolla en k etapas. Así, el experimento puede ser desarrollado en

$$n_1 \cdot n_2 \cdots n_k,$$

maneras.

2.4.2. Principio aditivo. Suponga que un experimento consta de dos etapas, tal que la primera se puede realizar de n_1 maneras y la segunda de n_2 maneras. Además considere que ambas etapas *no* se pueden realizar juntas. Entonces el experimento se puede llevar a cabo en

$$n_1 + n_2,$$

maneras.

OBSERVACIÓN 2.23. Evidentemente podemos extender lo anterior para un experimento que se desarrolla en k etapas. En cuyo caso tenemos que hay

$$n_1 + n_2 + \cdots + n_k,$$

formas de desarrollar el experimento.

DEFINICIÓN 2.24. Sea $A \neq \emptyset$ un conjunto finito, tal que $\#(A) = n$. Se denomina *variación simple* de orden k ($k \leq n$) a todo subconjunto de n elementos distinguiendo los elementos que lo componen y el orden en que son ordenados

EJEMPLO 2.25. Sea $A = \{\omega_1, \omega_2, \dots, \omega_n\}$. Para formar las variaciones de orden 2, digamos $(n)_2$ considere

$$\begin{pmatrix} \omega_1\omega_2 & \omega_1\omega_3 & \dots & \omega_1\omega_n \\ \omega_2\omega_1 & \omega_2\omega_3 & \dots & \omega_2\omega_n \\ \vdots & & & \\ \omega_n\omega_1 & \omega_n\omega_2 & \dots & \omega_n\omega_{n-1} \end{pmatrix},$$

note que esta matriz tiene n filas y $n-1$ columnas. Así, tenemos un total de $n(n-1)$ elementos, es decir

$$(n)_2 = n(n-1).$$

En general, el numero total de variaciones de orden k es dada por:

$$(n)_k = n(n-1) \cdots (n-k+1) = \frac{n!}{(n-k)!}.$$

OBSERVACIÓN 2.26. Recuerde que, para un entero positivo n , tenemos

$$n! = n(n-1) \cdots 2 \cdot 1 = n(n-1)!$$

Además $0! = 1$ por convención.

DEFINICIÓN 2.27. Sea $A \neq \emptyset$ un conjunto finito con $\#(A) = n$ se llama *variación con repetición* a toda variación en la que un mismo elemento puede aparecer n veces repetido.

Sea $A = \{\omega_1, \omega_2, \dots, \omega_n\}$, cada variación simple de orden $k-1$ da origen a n variaciones de orden k , luego se debe resolver el sistema de ecuaciones

$$\begin{aligned} (n)_1^* &= n \\ (n)_2^* &= n(n)_1^* \\ &\vdots \\ (n)_k^* &= n(n)_{k-1}^*. \end{aligned}$$

Es decir,

$$(n)_k^* = \underbrace{n \cdot n \cdots n}_{k \text{ veces}} = n^k.$$

EJEMPLO 2.28. Considere un grupo de k personas. ¿Cuántas listas posibles se pueden realizar con los días de sus cumpleaños? En efecto, tenemos

$$365 \cdot 365 \cdots 365 = 365^k.$$

EJEMPLO 2.29 (Problema del cumpleaños). Suponga un grupo de k personas. Se desea calcular la probabilidad de que 2 personas cumplan años el mismo día. Entonces, tenemos

$$\frac{(365)_k}{365^k} = \frac{365 \cdot 364 \cdots (365 - k + 1)}{365^k},$$

y por la regla del complemento,

$$P(\{\text{al menos un par de personas cumplan años el mismo día}\}) = 1 - \frac{(365)_k}{365^k}.$$

Es decir,

k	10	15	20	25	30	35	40	45	50
probabilidad	0.117	0.253	0.411	0.569	0.706	0.814	0.891	0.941	0.970

DEFINICIÓN 2.30 (Permutaciones). Sea $A \neq \emptyset$ un conjunto finito y $\#(A) = n$ se llama permutación a todo subconjunto de n elementos tales que las permutaciones se diferencian por el orden de sus elementos. Así,

$$P_{n,n} = (n)_n = n(n-1) \cdots 2 \cdot 1 = n!$$

OBSERVACIÓN 2.31. $P_{n,n}$ también se denominan permutaciones sin repetición

EJEMPLO 2.32. ¿De cuantas formas ubicar 5 personas en 5 sillas?

$$(5)_5 = 5! = 120$$

Cuando existen elementos iguales en A , las permutaciones se dicen *con* repetición. Considere,

$$A = \{\underbrace{\omega_1, \dots, \omega_1}_{n_1 \text{ veces}}, \underbrace{\omega_2, \dots, \omega_2}_{n_2 \text{ veces}}, \dots, \underbrace{\omega_k, \dots, \omega_k}_{n_k \text{ veces}}\},$$

tal que $\sum_{j=1}^k n_j = n$. Es decir, el número total de combinaciones es $n!$. Sea M el número de permutaciones distintas, entonces

$$M \prod_{j=1}^k n_j! = n! \quad \implies \quad M = \frac{n!}{\prod_{j=1}^k n_j!}.$$

EJEMPLO 2.33. Calcular el número de permutaciones distintas con las cifras $\{4, 7, 3, 4, 7, 7, 3\}$. De este modo,

$$M = \frac{7!}{2!2!3!} = \frac{5040}{2 \cdot 2 \cdot 6} = 210.$$

DEFINICIÓN 2.34 (Combinaciones). Sea $A \neq \emptyset$ y $\#(A) = n$, se denomina *combinación* de k elementos a todo conjunto diferente de k elementos tomado desde n .

OBSERVACIÓN 2.35. Dos combinaciones se distinguen entre sí por la naturaleza de los elementos que la componen sin interesar el *orden*.

EJEMPLO 2.36. Considere los objetos $\{a, b, c, d\}$ y $k = 2$. Entonces tenemos,

$$ab, ac, ad, bc, bd, cd,$$

es decir 6 elementos diferentes (note que ab y ba sólo difieren en el orden). Por tanto,

$$k! \binom{n}{k} = (n)_k \quad \implies \quad \binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Algunas propiedades del combinatorio son:

$$\begin{aligned} \binom{n}{k} &= \binom{n}{n-k} \\ \binom{n}{k} + \binom{n}{k+1} &= \binom{n+1}{k+1} \\ \binom{n}{k} &= \sum_{r=1}^{n-k-1} \binom{n-k}{k-1}. \end{aligned}$$

Además, estos coeficientes aparecen en múltiples ocasiones. Por ejemplo,

$$(a+b)^n = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k}.$$

2.5. Probabilidad condicional

EJEMPLO 2.37. Considere un lote con 80 artículos sin defectos y 20 defectuosos y suponga que se selecciona 2 artículos (a) con substitución, y (b) sin substitución. Defina los eventos:

$$A = \{\text{el 1er artículo es defectuoso}\},$$

$$B = \{\text{el 2do artículo es defectuoso}\}.$$

Cuando escogemos *con* substitución, tenemos:

$$P(A) = P(B) = \frac{20}{100} = \frac{1}{5}.$$

Cuando escogemos *sin* substitución, tenemos que:

$$P(A) = \frac{20}{100} = \frac{1}{5},$$

pero, ¿Cambia $P(B)$?

DEFINICIÓN 2.38. Si A y B son dos eventos en Ω y $P(B) > 0$, entonces la *probabilidad condicional* de A dado B , escrito $P(A|B)$ es

$$P(A|B) = \frac{P(A \cap B)}{P(B)}.$$

Note que $P(B|B) = 1$, es decir, B *actua* como Ω . En efecto, como $A = A \cap \Omega$, tenemos

$$P(A) = P(A|\Omega) = \frac{P(A \cap \Omega)}{P(\Omega)}.$$

La ocurrencia de A es calibrada con relación a B . En particular, si $A \cap B = \emptyset$, entonces

$$P(A|B) = P(B|A) = 0.$$

En el ejemplo anterior, se desea calcular $P(B|A) = 19/99$, pues si A ya ha ocurrido sólo quedan 19 defectuosos entre los 99 artículos.

Reexpresandola probabilidad condicional, tenemos

$$P(A \cap B) = P(A|B) P(B),$$

o bien

$$P(A \cap B) = P(B|A) P(A).$$

Las expresiones anteriores permiten “contornar” cálculos complicados, usando²

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}.$$

OBSERVACIÓN 2.39. El espacio de probabilidad definido por $\mathcal{F} \cap B$ permite notar que $P(A|B)$ es una función de probabilidad, es decir satisface:

- (a) $P(A|B) \geq 0$.
- (b) $P(\Omega|B) = 1$.
- (c) Para $\{A_n\}_{n \geq 1}$ sucesión disjunta

$$P\left(\bigcup_{n=1}^{\infty} A_n | B\right) = \sum_{n=1}^{\infty} P(A_n | B).$$

²Esto es un caso particular del Teorema de Bayes.

RESULTADO 2.40 (Probabilidad total). *Sea (Ω, \mathcal{F}, P) un espacio de probabilidad y sea C_1, C_2, \dots , una partición medible y contable tal que $P(C_i) \geq 0, \forall i$. Entonces, para todo $A \in \mathcal{F}$,*

$$P(A) = \sum_{i=1}^{\infty} P(A|C_i) P(C_i).$$

DEMOSTRACIÓN. Como los C_i 's forman una partición tenemos que

$$A = A \cap \Omega = A \cap \left(\bigcup_{i=1}^{\infty} C_i \right) = \bigcup_{i=1}^{\infty} (A \cap C_i).$$

Además,

$$P(A) = \sum_{i=1}^{\infty} P(A \cap C_i) = \sum_{i=1}^{\infty} P(A|C_i) P(C_i).$$

□

RESULTADO 2.41 (Teorema de Bayes). *Sea (Ω, \mathcal{F}, P) espacio de probabilidad y sea $\{C_i\}$ partición contable de Ω con $P(C_i) \geq 0, \forall i$. Entonces, para todo $A \in \mathcal{F}$, tenemos que*

$$P(C_i|A) = \frac{P(A|C_i) P(C_i)}{\sum_{k=1}^{\infty} P(A|C_k) P(C_k)},$$

siempre que $P(A) > 0$.

DEMOSTRACIÓN. Tenemos que

$$P(C_i|A) P(A) = P(A|C_i) P(C_i),$$

así

$$P(C_i|A) = \frac{P(A|C_i) P(C_i)}{P(A)}, \quad P(A) > 0,$$

desde el Teorema de probabilidad total, sigue que

$$P(C_i|A) = \frac{P(A|C_i) P(C_i)}{\sum_{k=1}^{\infty} P(A|C_k) P(C_k)}.$$

□

EJEMPLO 2.42. Considere un lote de 20 artículos defectuosos y 80 sin defectos, desde los que se escoge 2 artículos sin reemplazo. Sea

$$A = \{\text{el 1er artículo es defectuoso}\},$$

$$B = \{\text{el 2do artículo es defectuoso}\}.$$

Para calcular $P(B)$ podemos hacer

$$\begin{aligned} P(B) &= P(B|A) P(A) + P(B|A^c) P(A^c) = \frac{19}{99} \frac{1}{5} + \frac{20}{99} \frac{4}{5} \\ &= \frac{1}{5} \frac{1}{99} (19 + 20 \cdot 4) = \frac{1}{5}. \end{aligned}$$

2.6. Independencia estadística

DEFINICIÓN 2.43. Sea (Ω, \mathcal{F}, P) un espacio de probabilidad y sean $A, B \in \mathcal{F}$. Se dice que A y B son independientes si y sólo si

$$P(A \cap B) = P(A)P(B).$$

Naturalmente podemos entender la independencia del siguiente modo: “la ocurrencia de un evento B no tiene efecto en la probabilidad de otro evento A ”. Es decir,

$$P(A|B) = P(A).$$

Note también que,

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)} = \frac{P(A)P(B)}{P(B)} = P(B),$$

es decir, la ocurrencia de A no tiene efecto en B .

RESULTADO 2.44. Si A y B son independientes, entonces los siguientes pares también son independientes

- (a) A y B^c .
- (b) A^c y B .
- (c) A^c y B^c .

DEMOSTRACIÓN. Para probar (a), note que

$$\begin{aligned} P(A \cap B^c) &= P(A) - P(A \cap B) = P(A) - P(A)P(B) = P(A)(1 - P(B)) \\ &= P(A)P(B^c). \end{aligned}$$

Partes (b) y (c) son análogas y se dejan de ejercicio para el lector. \square

DEFINICIÓN 2.45. Una colección de eventos A_1, A_2, \dots, A_n es mutuamente independiente si para cualquier subcolección A_{i_1}, \dots, A_{i_k} , tenemos

$$P\left(\prod_{j=1}^k A_{i_j}\right) = \prod_{j=1}^k P(A_{i_j}).$$

EJEMPLO 2.46. Se lanzan 3 dados de distinto color: blanco, rojo y negro ¿Cuál es la probabilidad de que el dado blanco salga 3 y los otros dos no? Considere A , B y C los eventos

A : resultado del dado blanco es 3,

B : resultado del dado rojo es 3,

C : resultado del dado negro es 3,

tenemos $P(A) = P(B) = P(C) = 1/6$ y se pide calcular

$$P(A \cap B^c \cap C^c) = P(A)P(B^c)P(C^c) = \frac{1}{6} \frac{5}{6} \frac{5}{6} = \frac{25}{216}.$$

Bibliografía

- Anscombe, F.J. (1973) Graphs in statistical analysis. *The American Statistician* **27**, 17-21.
- Barlow, J.L. (1993). Numerical aspects of solving linear least squares problems. En: *Handbook of Statistics, Vol. 9*, C.R. Rao (Ed.). Elsevier, pp. 303-373.
- Bolfarine, H., Sandoval, M.C. (2001). *Introdução à Inferência Estatística*. Sociedade Brasileira de Matemática, Rio de Janeiro.
- Casella, G., Berger, R.L. (2002). *Statistical Inference (2nd Ed.)*. Duxbury, Pacific Grove.
- Chan, T.F., Lewis, J.G. (1979). Computing standard deviations: Accuracy. *Communications of the Association for Computing Machinery* **22**, 526-531.
- Chan, T.F., Golub, G.H., LeVeque, R.J. (1983). Algorithms for computing the sample variance: Analysis and recommendations. *The American Statistician* **37**, 242-247.
- Clarke, M.R.B. (1971). Algorithm AS 41: Updating the sample mean and dispersion matrix. *Applied Statistics* **20**, 206-209.
- Cleveland, W.S. (1993). *Visualizing Data*. Hobart Press, New Jersey.
- Frery, A.C., Cribari-Neto, F. (2005). *Elementos de Estatística Computacional Usando Plataformas de Software Livre/Gratuito*. IMPA, Rio de Janeiro.
- Fisher, R.A. (1936). The use of multiple measurements in taxonomic problems. *Annals of Eugenics* **7**, 179-188.
- Galbiati, J. (2012). *Tablas de Probabilidad (8a Ed.)*. Instituto de Estadística, Pontificia Universidad Católica de Valparaíso, Chile.
- Graillat, S. (2009). Accurate floating-point product and exponentiation. *IEEE Transactions on Computers* **58**, 994-1000.
- Grossman, S.I., Turner, J.E. (1974). *Mathematics for the Biological Sciences*. MacMillan Publishing, New York.
- Jambu, M. (1991). *Exploratory and Multivariate Data Analysis*. Academic Press, Boston.
- Keeling, K.B., Pavur, R.J. (2007). A comparative study of the reliability of nine statistical software packages. *Computational Statistics & Data Analysis* **51**, 3811-3831.
- Knuth, D.E. (1997). *The Art of Computer Programming: Vol. 1, Fundamental Algorithms*. Addison-Wesley, Massachusetts.
- Mardia, K.V. (1970). Measures of multivariate skewness and kurtosis with applications. *Biometrika* **57**, 519-530.
- Mardia, K.V. (1974). Application of some measures of multivariate skewness and kurtosis in testing normality and robustness studies. *Sankhyā, Series B* **36**, 115-128.

- McCullough, B.D., Wilson, B. (1999). On the accuracy of statistical procedures in Microsoft Excel 97. *Computational Statistics & Data Analysis* **31**, 27-37.
- McCullough, B.D., Wilson, B. (2002). On the accuracy of statistical procedures in Microsoft Excel 2000 and Excel XP. *Computational Statistics & Data Analysis* **40**, 713-721.
- McCullough, B.D., Wilson, B. (2005). On the accuracy of statistical procedures in Microsoft Excel 2003. *Computational Statistics & Data Analysis* **49**, 1244-1252.
- Mood, A.M., Graybill, F.A., Boes, D.C. (1974). *Introduction to the Theory of Statistics (3rd Ed.)*. McGraw-Hill, New York.
- Osorio, F., Ogeda, A. (2022). *fastmatrix*: Fast computation of some matrices useful in statistics. R package version 0.4. URL: faosorios.github.io/fastmatrix/
- Panaretos, V.M. (2016). *Statistics for Mathematicians: A Rigorous First Course*. Birkhäuser, Laussane.
- Pébay, P., Terriberry, T.B., Kolla, H., Bennett, J. (2016). Numerically stable, scalable formulas for parallel and online computation of higher order multivariate central moments with arbitrary weights. *Computational Statistics* **31**, 1305-1325.
- Press, W.H., Teukolsky, S.A., Vetterling, W.T., Flannery, B.P. (1992). *Numerical Recipes in C: The Art of Scientific Computing*. Cambridge University Press, Cambridge.
- R Core Team (2019). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL: www.R-project.org
- Spicer, C.C. (1972). Algorithm AS 52: Calculation of power sums of deviations about the mean. *Applied Statistics* **21**, 226-227.
- Thisted, R.A. (1988). *Elements of Statistical Computing*. Chapman & Hall, Boca Raton.
- Tukey, J.W. (1977). *Exploratory Data Analysis*. Addison-Wesley, Massachusetts.
- Venables, W.N., Ripley, B.D. (1994). *Modern Applied Statistics with S-Plus*. Springer, New York.
- Wasserman, L. (2004). *All of Statistics: A Concise Course in Statistical Inference*. Springer, New York.
- Welford, B.P. (1962). Note on a method for calculating corrected sums of squares and products. *Technometrics* **4**, 419-420.
- West, D.H.D. (1979). Updating mean and variances estimates: An improved method. *Communications of the Association for Computing Machinery* **22**, 532-535.
- Wilkinson, L. (1999). *The Grammar of Graphics*. Springer, New York.