

BIOESTADÍSTICA APLICADA CON R Y RKTTEACHING

Santiago Angulo Díaz-Parreño (sangulo@ceu.es)
Euardo López Ramírez (elopez@ceu.es)
José Rojo Montijano (jrojo.eps@ceu.es)
Anselmo Romero Limón (arlimon@ceu.es)
Alfredo Sánchez Alberca (asalber@ceu.es)
Susana Victoria Rodríguez (victoria.eps@ceu.es)

Septiembre de 2014

Bioestadística Aplicada con R y RKTeaching

Alfredo Sánchez Alberca (asalber@gmail.com).



Esta obra está bajo una licencia Reconocimiento – No comercial – Compartir bajo la misma licencia 3.0 España de Creative Commons. Para ver una copia de esta licencia, visite <http://creativecommons.org/licenses/by-nc-sa/3.0/es/>.

Con esta licencia eres libre de:

- Copiar, distribuir y mostrar este trabajo.
- Realizar modificaciones de este trabajo.

Bajo las siguientes condiciones:



Reconocimiento. Debe reconocer los créditos de la obra de la manera especificada por el autor o el licenciador (pero no de una manera que sugiera que tiene su apoyo o apoyan el uso que hace de su obra).



No comercial. No puede utilizar esta obra para fines comerciales.



Compartir bajo la misma licencia. Si altera o transforma esta obra, o genera una obra derivada, sólo puede distribuir la obra generada bajo una licencia idéntica a ésta.

- Al reutilizar o distribuir la obra, tiene que dejar bien claro los términos de la licencia de esta obra.
 - Alguna de estas condiciones puede no aplicarse si se obtiene el permiso del titular de los derechos de autor
 - Nada en esta licencia menoscaba o restringe los derechos morales del autor.
-

Índice general

1. Introducción a R y RKWard	1
1.1. Introducción	1
1.2. Instalación	2
1.2.1. Instalación de R	2
1.2.2. Instalación de la interfaz gráfica RKWard y el paquete rkTeaching	2
1.3. Arranque	3
1.4. Tipos de datos y operadores aritméticos y lógicos	4
1.5. Introducción y manipulación de datos	5
1.5.1. Introducción de datos en línea de comandos	5
1.5.2. Introducción de datos en RKWard	6
1.5.3. Ponderación de datos	7
1.5.4. Guardar datos	7
1.5.5. Abrir datos	8
1.5.6. Eliminación de datos	9
1.6. Transformación de datos	9
1.6.1. Filtrado de datos	9
1.6.2. Cálculo de variables	9
1.6.3. Recodificación de variables	10
1.7. Manipulación de ficheros de resultados	10
1.7.1. Guardar los resultados	11
1.7.2. Limpiar la ventana de resultados	11
1.8. Manipulación de guiones de comandos	11
1.8.1. Creación de un guión de comandos	11
1.8.2. Guardar un guión de comandos	12
1.8.3. Abrir un guión de comandos	12
1.9. Ayuda	12
1.10. Ejercicios resueltos	13
1.11. Ejercicios propuestos	15
2. Distribuciones de Frecuencias y Representaciones Gráficas	17
2.1. Fundamentos teóricos	17
2.1.1. Cálculo de Frecuencias	17
2.1.2. Representaciones Gráficas	19
2.2. Ejercicios resueltos	23
2.3. Ejercicios propuestos	25
3. Estadísticos Muestrales	27
3.1. Fundamentos teóricos	27
3.1.1. Medidas de posición	27
3.1.2. Medidas de dispersión	28
3.1.3. Medidas de forma	29

3.1.4. Estadísticos de variables en las que se definen grupos	30
3.2. Ejercicios resueltos	31
3.3. Ejercicios propuestos	32
4. Regresión Lineal Simple y Correlación	35
4.1. Fundamentos teóricos	35
4.1.1. Regresión	35
4.1.2. Correlación	39
4.2. Ejercicios resueltos	42
4.3. Ejercicios propuestos	46
5. Regresión no lineal	49
5.1. Fundamentos teóricos	49
5.2. Ejercicios resueltos	51
5.3. Ejercicios propuestos	54
6. Probabilidad	55
6.1. Fundamentos teóricos	55
6.1.1. Introducción	55
6.1.2. Experimentos y sucesos aleatorios	55
6.1.3. Definición de probabilidad	58
6.1.4. Probabilidad condicionada	60
6.1.5. Espacios probabilísticos	61
6.1.6. Teorema de la probabilidad total	62
6.1.7. Teorema de Bayes	63
6.1.8. Tests diagnósticos	64
6.2. Ejercicios resueltos	66
6.3. Ejercicios propuestos	71
7. Variables Aleatorias Discretas	73
7.1. Fundamentos teóricos	73
7.1.1. Variables Aleatorias	73
7.1.2. Variables Aleatorias Discretas (v.a.d.)	73
7.2. Ejercicios resueltos	77
7.3. Ejercicios propuestos	80
8. Variables Aleatorias Continuas	81
8.1. Fundamentos teóricos	81
8.1.1. Variables Aleatorias	81
8.1.2. Variables Aleatorias Continuas (v.a.c.)	81
8.2. Ejercicios resueltos	87
8.3. Ejercicios propuestos	92
9. Intervalos de Confianza para Medias y Proporciones	93
9.1. Fundamentos teóricos	93
9.1.1. Inferencia Estadística y Estimación de Parámetros	93
9.1.2. Intervalos de Confianza	93
9.2. Ejercicios resueltos	98
9.3. Ejercicios propuestos	100
10. Intervalos de Confianza para la Comparación de 2 Poblaciones	103
10.1. Fundamentos teóricos	103
10.1.1. Inferencia Estadística y Estimación de Parámetros	103
10.1.2. Intervalos de Confianza	103
10.2. Ejercicios resueltos	109
10.3. Ejercicios propuestos	111

Introducción a R y RKWard

1 Introducción

La gran potencia de cálculo alcanzada por los ordenadores ha convertido a los mismos en poderosas herramientas al servicio de todas aquellas disciplinas que, como la estadística, requieren manejar un gran volumen de datos. Actualmente, prácticamente nadie se plantea hacer un estudio estadístico serio sin la ayuda de un buen programa de análisis estadístico.

R es un potente lenguaje de programación que incluye multitud de funciones para la representación el análisis de datos. Fue desarrollado por Robert Gentleman y Ross Ihaka en la Universidad de Auckland en Nueva Zelanda, aunque actualmente es mantenido por una enorme comunidad científica en todo el mundo.



Las ventajas de R frente a otros programas habituales de análisis de datos, como pueden ser SPSS, SAS, SPlus, Matlab o Minitab, son múltiples:

- Es software libre y por tanto gratuito. Puede descargarse desde la web <http://www.r-project.org/>.
- Es multiplataforma. Existen versiones para Windows, Macintosh, Linux y otras plataformas.
- Está avalado y en constante desarrollo por una amplia comunidad científica que lo utiliza como estándar para el análisis de datos.
- Cuenta con multitud de paquetes para todo tipo de análisis estadísticos y representaciones gráficas, desde los más habituales, hasta los más novedosos y sofisticados que no incluyen otros programas. Los paquetes están organizados y documentados en un repositorio CRAN (Comprehensive R Archive Network) desde donde pueden descargarse libremente. En España hay una copia de este repositorio en la web <http://cran.es.r-project.org/>.
- Es programable, lo que permite que el usuario pueda crear fácilmente sus propias funciones o paquetes para análisis de datos específicos.
- Existen multitud de libros, manuales y tutoriales libres que permiten su aprendizaje e ilustran el análisis estadístico de datos en distintas disciplinas científicas como las matemáticas, la física, la biología, la psicología, la medicina, etc.

Por defecto el entorno de trabajo de R es en línea de comandos, lo que significa que los cálculos y los análisis se realizan mediante comandos o instrucciones que el usuario teclea en una ventana de texto. No obstante, existen distintas interfaces gráficas de usuario que facilitan su uso, sobre todo para usuarios noveles. La interfaz gráfica que se utilizará para realizar estas prácticas será *RKward*, desarrollada por Thomas Friedrichsmeier, junto al paquete *rkTeaching* especialmente desarrollado por el departamento de Matemáticas de la Universidad San Pablo CEU para la docencia de estadística.

El objetivo de esta práctica es introducir al alumno en la utilización de este programa, enseñándole a realizar las operaciones básicas más habituales de carga y manipulación de datos.

2 Instalación

2.1 Instalación de R

Linux En la distribución Debian y cualquiera de sus derivadas (Ubuntu, Kubuntu, etc.) basta con teclear en la línea de comandos

```
> sudo apt-get install r-base-html r-cran-rcmdr r-cran-rodbc r-doc-html r-recommended
```

Windows Descargar de <http://cran.es.r-project.org/bin/windows/base/release.htm> el programa de instalación de R, ejecutarlo y seguir las instrucciones de instalación.

2.2 Instalación de la interfaz gráfica RKward y el paquete rkTeaching

La interfaz gráfica de usuario RKward puede descargarse desde la web <http://rkward.sourceforge.net/> donde se indican las instrucciones para instalarlo en cada plataforma.

Para Windows se recomienda seleccionar el paquete de instalación completa que incorpora R, las librerías gráficas de KDE y el propio RKward.

R dispone de una gran librería de paquetes que incorporan nuevas funciones y procedimientos. En la instalación base de R vienen ya cargados los procedimientos y funciones para los análisis más comunes, pero en ocasiones, para otros análisis será necesario cargar algún paquete adicional como por ejemplo el paquete *rkTeaching* que incorpora un nuevo menú a RKward con la mayoría de los análisis que se realizarán en estas prácticas.

Para instalar el paquete *rkTeaching*, basta con descargarlo desde la dirección http://asalber.github.io/rkTeaching_es/, arrancar R o RKward y, en la consola de comandos, teclear el comando

```
> setwd("ruta_a_descargas")
> install.packages("rkTeaching", repos=NULL, dep=TRUE)
```

La instalación de cualquier otro paquete se realiza con el mismo comando, cambiando el nombre del paquete por el deseado.

En RKward, también puede instalarse desde la ventana de R mediante el menú **Preferencias** ▶ **Configurar paquetes**. Con esto aparecerá una ventana donde se muestran los paquetes instalados localmente. Para cargar un paquete instalado localmente basta con seleccionarlo y hacer clic sobre el botón **Cargar**. En esa misma ventana aparece una solapa **Install/Update/Remove** que permite instalar nuevos paquetes desde un repositorio de R. Al hacer clic sobre esta solapa se abrirá una conexión a internet y aparecerá una ventana con los distintos repositorios disponibles. Normalmente seleccionaremos en más cercano geográficamente, en nuestro caso **Spain(Madrid)**. Después aparecerá una lista de paquetes instalados y nuevos. Para instalar un paquete nuevo basta con seleccionarlo y hacer clic en el botón **Aceptar**. Una vez instalado localmente, podrá cargarse como se ha indicado antes.

3 Arranque

Como cualquier otra aplicación de Windows, para arrancar el programa hay que hacer clic sobre la opción correspondiente del menú Inicio►Programas►RKWard, o bien sobre el icono de escritorio



Al arrancar, aparece la ventana de bienvenida de RKWard (figura 1.1).

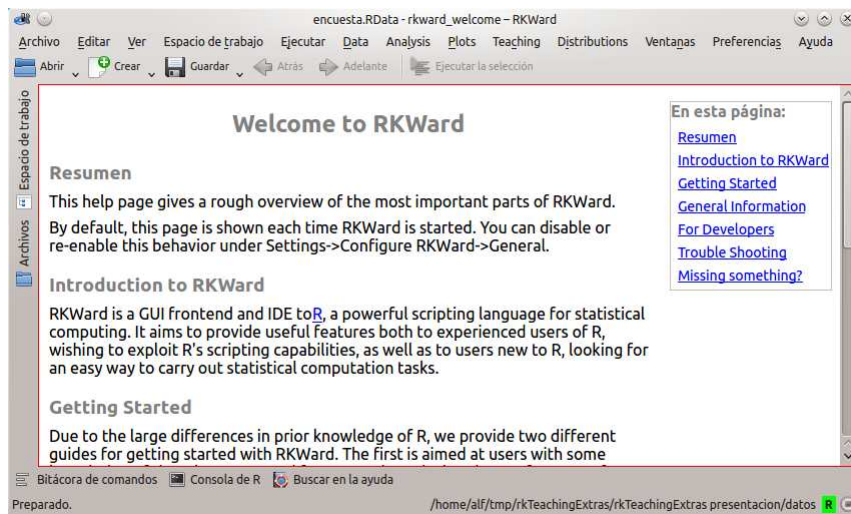


Figura 1.1 – Interfaz gráfica de usuario de RKWard.

La interfaz gráfica de usuario RKWard consta de los siguientes elementos:

- **Barra de menús.** Contiene distintos menús con operaciones que pueden realizarse con R. Si se ha instalado el paquete rkTeaching debe de aparecer el menú Teaching.
- **Barra de botones.** Contiene botones para abrir, crear y guardar conjuntos de datos, espacios de trabajo y guiones de comandos.
- **Ventana principal.** Es la ventana central donde aparecerán la ventana de introducción de datos, los resultados de los comandos ejecutados o de las búsquedas realizadas.
- **Espacio de trabajo.** Es una ventana desplegable al hacer clic sobre la solapa situada en el lado izquierdo que contiene todos los elementos del espacio de trabajo de R. Entre estos elementos aparecen los paquetes cargados, los conjuntos de datos y las variables que contienen los datos de la sesión actual.
- **Bitácora de comandos** Es una solapa desplegable situada en la parte inferior donde aparece un registro de todas las acciones realizadas o comandos ejecutados en la sesión de trabajo actual. Cada vez que se seleccione un menú que lleve asociado la ejecución de algún comando, dicho comando aparecerá en esta ventana. Esto permite modificar fácilmente los parámetros del comando y volver a ejecutarlo rápidamente sin necesidad de volver al menú.
- **Consola de R** Es una solapa desplegable situada también en la parte inferior que da acceso al intérprete de comandos de R. En esta ventana pueden teclearse y ejecutarse directamente los comandos de R.

- **Buscar en la ayuda** Es una solapa desplegable situada en la parte inferior que permite hacer búsquedas sobre comandos de R o de algún paquete.
- **Mensajes.** Es la línea de texto que aparece en la parte inferior, donde se muestra información adicional sobre errores, advertencias u otra información auxiliar al ejecutar un comando, así como la ruta del espacio de trabajo activo.

4 Tipos de datos y operadores aritméticos y lógicos

En R existen distintos tipos de datos. Los más básicos son:

Numeric : Es cualquier número decimal. Se utiliza el punto como separador de decimales. Por defecto, cualquier número que se teclee tomará este tipo.

Integer : Es cualquier número entero. Para convertir un número de tipo Numeric en un entero se utiliza el comando `as.integer()`

Logical : Puede tomar cualquiera de los dos valores lógicos TRUE (verdadero) o FALSE (falso).

Character : Es cualquier cadena de caracteres alfanuméricos. Deben introducirse entre comillas. Para convertir cualquier número en una cadena de caracteres se utiliza el comando `as.character()`.

Los valores de estos tipos de datos pueden operarse utilizando distintos operadores o funciones predefinidas para cada tipo de datos. Los más habituales son:

Operadores aritméticos : + (suma), - (resta), * (producto), / (cociente), ^ (potencia).

Operadores de comparación : > (mayor), < (menor), >= (mayor o igual), <= (menor o igual), == (igual), != (distinto).

Operadores lógicos : & (conjunción y), | (disyunción o), ! (negación no).

Funciones predefinidas : `sqrt()` (raíz cuadrada), `abs()` (valor absoluto), `log()` (logaritmo neperiano), `exp()` (exponencial), `sin()` (seno), `cos()` (coseno), `tan()` (tangente).

Al evaluar las expresiones aritméticas existe un orden de prioridad entre los operadores de manera que primero se evalúan las funciones predefinidas, luego las potencias, luego los productos y cocientes, luego las sumas y restas, luego los operadores de comparación, luego las negaciones, luego las conjunciones y finalmente las disyunciones. Para forzar un orden de evaluación distinto del predefinido se pueden usar paréntesis. Por ejemplo

```
> 2^2+4/2
[1] 6
> (2^2+4)/2
[1] 4
> 2^(2+4/2)
[1] 16
> 2^(2+4)/2
[1] 32
> 2^((2+4)/2)
[1] 8
```

También es posible asignar valores a variables mediante el operador de asignación =. Una vez definidas, las variables pueden usarse en cualquier expresión aritmética o lógica. Por ejemplo,

```
> x=2
> y=x+2
> y
[1] 4
```



```
> y>x
[1] TRUE
> x>=y
[1] FALSE
> x==y-2
[1] TRUE
> x!=0 & !y<x
[1] TRUE
```

5 Introducción y manipulación de datos

Antes de realizar cualquier análisis de datos hay que introducir los datos que se quieren analizar.

5.1 Introducción de datos en línea de comandos

Existen muchas formas de introducir datos en R pero aquí sólo veremos las más habituales. La forma más rápida de introducir datos es usar la consola de R para crear un vector de datos mediante el comando `c()`. Por ejemplo, para introducir las notas de 5 alumnos se debe teclear en la consola de R

```
> nota = c(5.6, 7.2, 3.5, 8.1, 6.4)
```

Esto crea el vector `nota` con el que posteriormente se pueden realizar cálculos como por ejemplo la media

```
> mean(nota)
[1] 6.16
```

Otra forma habitual de introducir los datos de una muestra es crear un conjunto de datos mediante el comando `data.frame()`. Por ejemplo, para crear un conjunto de datos a partir de las notas anteriores, hay que teclear

```
> curso = data.frame(nota)
```

Esto crea una matriz de datos en la que cada columna se corresponde con una variable y cada fila con un individuo de la muestra. En el ejemplo la matriz `curso` sólo tendría una columna que se correspondería con las notas y 5 filas, cada una de ellas correspondiente a un alumno de la muestra. Es posible acceder a las variables de un conjunto de datos con el operador dolar `$`. Por ejemplo, para acceder a las notas hay que teclear

```
> curso$nota
[1] 5.6 7.2 3.5 8.1 6.4
```

Es fácil añadir nuevas variables a un conjunto de datos, pero siempre deben tener el mismo tamaño muestral. Por ejemplo, para añadir una nueva variable con el grupo (mañana o tarde) de los alumnos, hay que teclear

```
> curso$grupo = c("m", "t", "t", "m", "m")
```

Ahora el conjunto de datos `curso` tendría dos columnas, una para la nota y otra para el grupo de los alumnos. Tecleando el nombre de cualquier objeto, se muestra su información:

```
> curso
  nota grupo
1  5.6     m
2  7.2     t
3  3.5     t
4  8.1     m
5  6.4     m
```

Cuando se introducen datos se puede utilizar el código NA (not available), para indicar la ausencia del dato.

Las variables definidas en cada sesión de trabajo quedan almacenadas en la memoria interna de R en lo que se conoce como *espacio de trabajo*. Es posible obtener un listado de todos los objetos almacenados en el espacio de trabajo mediante los comandos `ls()`. Si se desea más información, el comando `ls.str()` además de mostrar los objetos de la memoria indica sus tipos y sus valores.

```
> ls()
[1] "curso" "nota"  "x"      "y"
> ls.str()
curso : 'data.frame': 5 obs. of 2 variables:
 $ nota : num 5.6 7.2 3.5 8.1 6.4
 $ grupo: chr " m " " t " " t " " m " ...
nota : num [1:5] 5.6 7.2 3.5 8.1 6.4
x : num 2
y : num 4
```

Para eliminar un objeto de la memoria se utiliza el comando `rm()`.

```
> ls()
[1] "curso" "nota"  "x"      "y"
> rm(x,y)
> ls()
[1] "curso" "nota"
```

5.2 Introducción de datos en RKWard

RKWard dispone de una interfaz gráfica para introducir los datos sin necesidad de saberse los comandos anteriores. Para ello hay que ir al menú Archivo ▶ Nuevo ▶ Conjunto de datos. Con esto aparecerá una ventana donde hay que darle un nombre al conjunto de datos y tras esto aparece la ventana de la figura 1.2 con una tabla en la que se pueden introducir los datos de la muestra. Al igual que antes, cada variable debe introducirse en una columna y cada individuo en una fila.

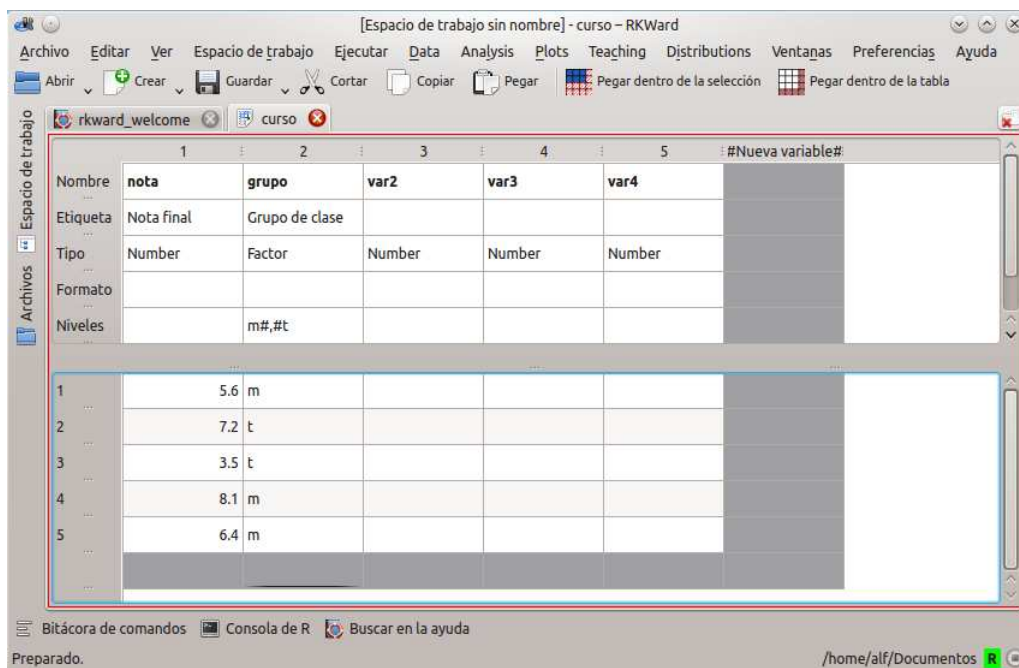


Figura 1.2 – Ventana de introducción de datos

Haciendo clic en las casillas de la cabecera cada fila es posible cambiar el nombre de la variable, ponerle una etiqueta, su tipo, su formato y los niveles en caso de tratarse de un factor o variable categórica. Los nombres de variables deben comenzar con una letra o un punto y pueden contener cualquier letra, punto, subrayado (_) o número. En particular, no se pueden utilizar espacios en blanco. Además, R distingue entre mayúsculas y minúsculas.

Una vez definida la variable, para introducir los datos basta con teclearlos en las casillas que aparecen más abajo en la misma columna.

R permite definir más de un conjunto de datos en un mismo espacio de trabajo.

Los objetos definidos en el espacio de trabajo pueden verse haciendo clic en la solapa Espacio de trabajo. Para editar una variable o un conjunto de datos basta con hacer doble clic sobre él. También puede obtenerse un resumen como el que se muestra en la figura 1.3 haciendo clic en el botón derecho y seleccionando ver en el menú contextual que aparece.

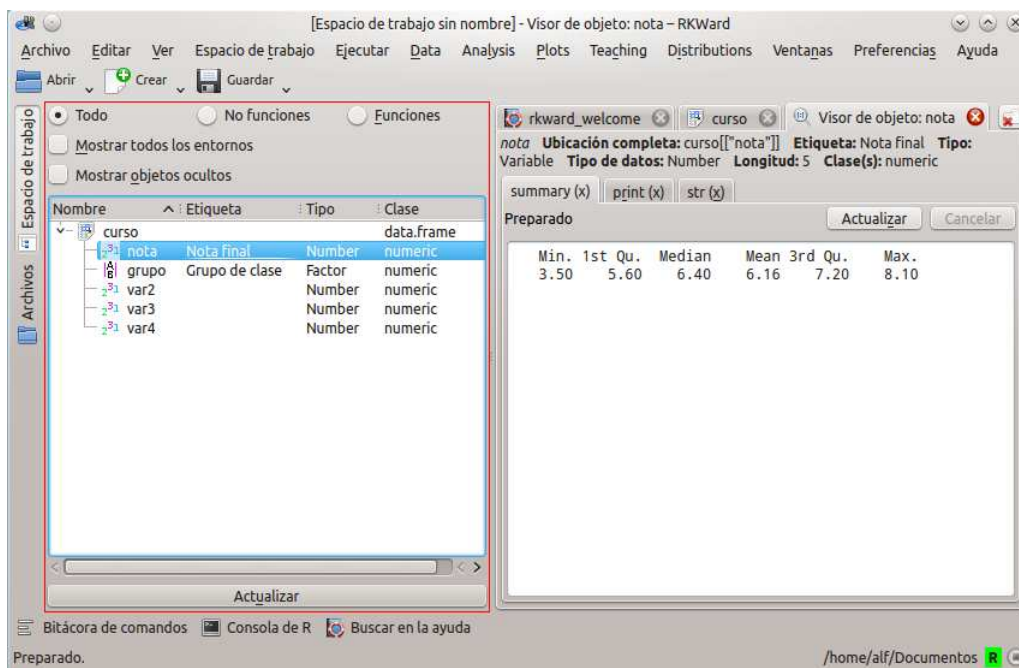


Figura 1.3 – Ventana de resumen descriptivo de un conjunto de datos

5.3 Ponderación de datos

Cuando una variable o un conjunto de datos tiene unos pocos valores que se repiten mucho, en lugar de teclearlos es más rápido indicar los valores y ponderarlos por sus frecuencias. Para ello se utiliza el menú Teaching ▶ Datos ▶ Ponerar datos. Al seleccionarlo aparece una ventana donde hay que seleccionar el conjunto de datos a ponderar, la variable numérica de dicho conjunto de datos que contiene las frecuencias de ponderación, e indicar un nombre para el nuevo conjunto de datos. Por ejemplo, si en una clase hay 20 chicas y 30 chicos, se puede crear un conjunto de datos con la variables sexo y frecuencia, tal y como se muestra en la figura 1.4, y después llamar al menú de ponderación con los datos que aparecen la figura 1.5.

5.4 Guardar datos

Una vez introducidos los datos, conviene guardarlos en un fichero para no tener que volver a introducirlos en futuras sesiones. Para guardar los conjunto de datos definidos en el espacio de trabajo, se utiliza el menú Espacio de trabajo ▶ Guardar espacio de trabajo. Con esto aparece una ventana

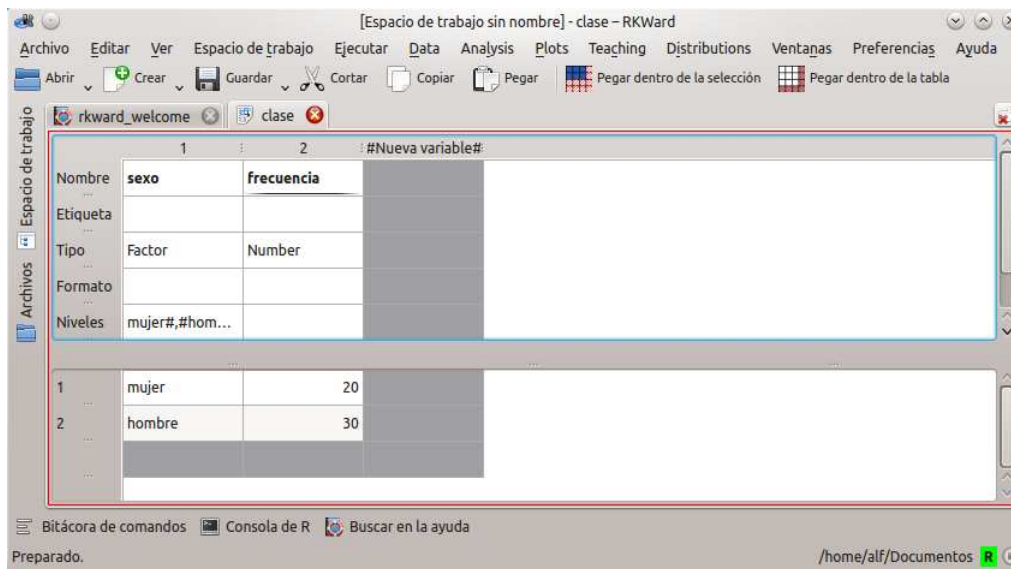


Figura 1.4 – Conjunto de datos preparado para ser ponderado

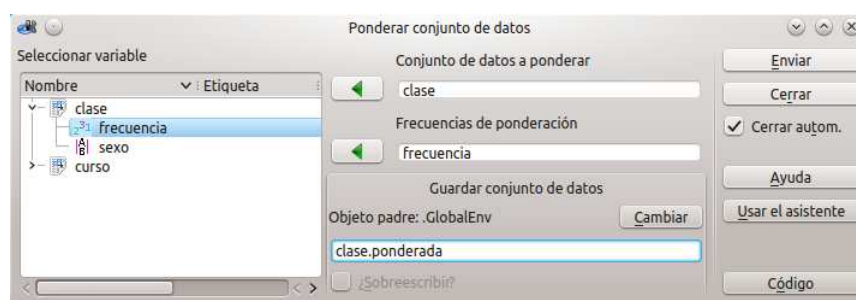


Figura 1.5 – Ventana de ponderación de datos

donde hay que darle un nombre al fichero y seleccionar la carpeta donde se guardará. Los conjuntos de datos se guardan siempre en ficheros de R con extensión `rda` o `rData`.

También es posible guardar los datos en un fichero de texto plano mediante el menú **Archivo** ▶ **Exportar** ▶ **Export tabular data**. Tras esto aparece una ventana donde hay que seleccionar el conjunto de datos a exportar, darle un nombre al fichero de texto y seleccionar la carpeta donde se guardará. Esta ventana contiene también solapas donde se puede indicar entre otras cosas si incluir los nombres de las variables o no, el separador de decimales o el separador de los datos, que puede ser un espacio, tabuladores, comas u otro carácter.

5.5 Abrir datos

Si los datos con los que se pretende trabajar ya están guardados en un fichero de R, entonces tendremos que abrir dicho fichero. Para ello se utiliza el **Espacio de trabajo** ▶ **Abrir espacio de trabajo** y en la ventana que aparece se selecciona el fichero que se desea abrir. Automáticamente se cargará el conjunto de datos del fichero y pasará a ser el conjunto de datos activo.

También es posible cargar datos de ficheros con otros formatos, como por ejemplo un fichero de texto. Para ello se utiliza el menú **Archivo** ▶ **Importar** ▶ **Importar datos** y en la ventana que aparece se selecciona el fichero de texto que se desea abrir y en el cuadro desplegable del formato de archivo se debes seleccionar **Text**. Después aparecerá una ventana donde habrá que darle un nombre al conjunto de datos y seleccionar el tipo de separador y si los nombres de las variables aparecen en la primera línea del fichero.

5.6 Eliminación de datos

Para eliminar una variable del conjunto de datos primero hay que editar el conjunto de datos, y después, en la ventana de edición de datos, hay que hacer clic con el botón derecho del ratón sobre la cabecera de la columna correspondiente y seleccionar en el menú contextual que aparece **Borrar** esta variable.

Para eliminar individuos del conjunto de datos que hacer clic con el botón derecho del ratón sobre la cabecera de la fila correspondiente y seleccionar en el menú contextual que aparece **Borrar** esta fila.

En la ventana del espacio de trabajo también es posible borrar cualquier objeto del espacio de trabajo de R haciendo clic con el botón derecho del ratón sobre él y seleccionando el menú **Eliminar**.

6 Transformación de datos

A menudo en los análisis hay que realizar transformaciones en los datos originales. A continuación se presentan las transformaciones más habituales.

6.1 Filtrado de datos

Cuando se desea realizar un análisis con un subconjunto de individuos del conjunto de datos activo que cumplen una determinada condición es posible filtrar el conjunto de datos para quedarse con esos individuos. Para ello se utiliza el menú **Teaching ▸ Datos ▸ Filtrar**. Con esto aparece un cuadro de diálogo en el que hay que seleccionar el conjunto de datos que se desea filtrar, y en el cuadro de texto **Condición de selección** indicar la condición lógica que tienen que cumplir los individuos seleccionados. También hay que indicar el nombre del nuevo conjunto de datos. Por ejemplo, para seleccionar los alumnos del grupo de la mañana habría que indicar la condición `grupo=="m"` tal y como se muestra en la figura 1.6.

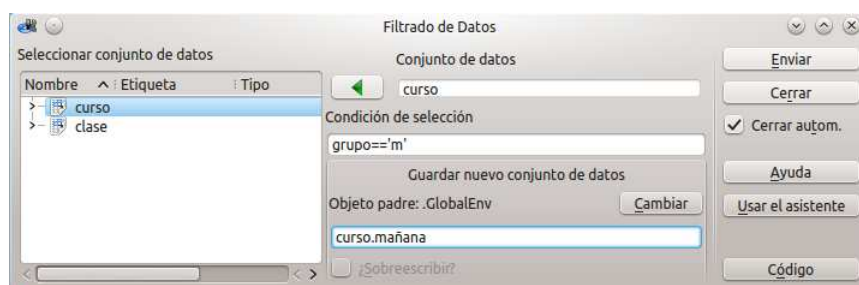


Figura 1.6 – Ventana de filtrado de datos.

6.2 Cálculo de variables

Para calcular una nueva variable a partir de otras ya existentes en el espacio de trabajo de R se utiliza el menú **Teaching ▸ Datos ▸ Calcular variable**. Con esto aparece un cuadro de diálogo en el que hay que introducir la expresión a partir de la que se calculará la nueva variable en el cuadro de texto **Expresión de cálculo**, e indicar el nombre de la nueva variable. La expresión de cálculo puede ser cualquier expresión aritmética o lógica de R, en las que pueden utilizarse cualquiera de las variables del espacio de trabajo de R. Por ejemplo, para eliminar los decimales de la variable `nota` podría crearse una nueva variable `puntuacion` multiplicando por 10 las notas, tal y como se muestra en la figura 1.7.

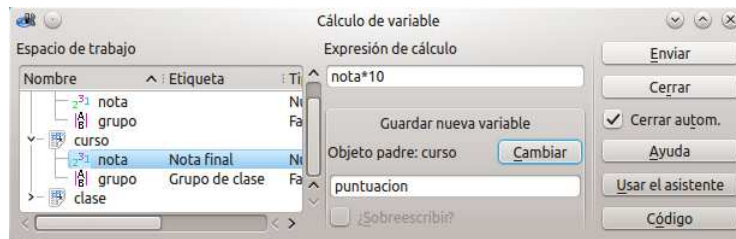


Figura 1.7 – Ventana de cálculo de nuevas variables.

6.3 Recodificación de variables

Otra transformación habitual es la recodificación de variables que permite transformar los valores de una variable de acuerdo a un conjunto de reglas de reescritura. Normalmente se utiliza para convertir una variable numérica en una variable categórica que pueda usarse como un factor.

Para recodificar una variable se utiliza el menú Teaching ▶ Datos ▶ Recodificar variable. Con esto aparece una ventana en la que hay que seleccionar la variable que se desea recodificar, indicar el nombre de la nueva variable recodificada e introducir las reglas de recodificación en el cuadro de texto Reglas de recodificación. Las reglas de recodificación siempre siguen la sintaxis **valor o rango de valores = nuevo valor** y pueden introducirse tantas reglas como se desee, cada una en una línea. Al lado izquierdo de la igualdad puede introducirse un único valor, varios valores separados por comas, o un rango de valores indicando el límite inferior y el límite superior del intervalo separados por el operador **:**. A la hora de definir el límite inferior puede utilizarse la palabra clave **lo** para referirse al menor de los valores de la muestra y **hi** para referirse al mayor de los valores. Por ejemplo, para recodificar la variable **nota** en categorías correspondientes a las calificaciones ([0,5) Suspenso, [5,7) Aprobado, [7,9) Notable y [9,10] Sobresaliente), habría que introducir las reglas que se muestran en la figura 1.8. Después, en la ventana de introducción de datos, se pueden renombrar los niveles del factor introduciendo el valor suspenso para la categoría 1, aprobado para la categoría 2, notable para la categoría 3 y sobresaliente para la categoría 4.

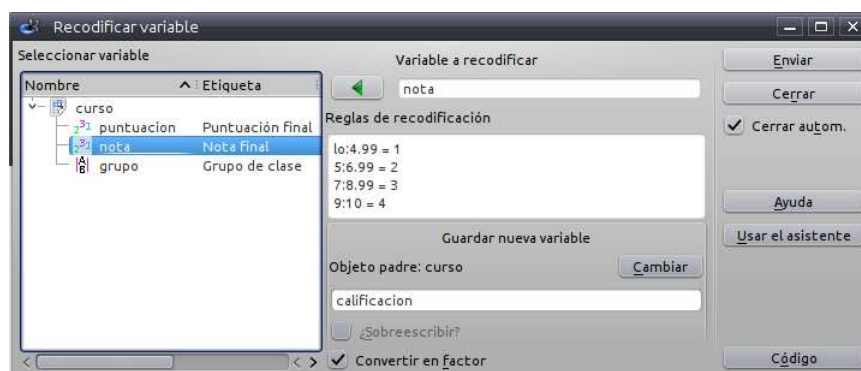


Figura 1.8 – Ventana de recodificación de variables

7 Manipulación de ficheros de resultados

7.1 Guardar los resultados

Cada vez que se ejecuta un comando de R, bien en la consola de comandos o a través de un menú, el comando ejecutado y su salida quedan registrados en la bitácora de comandos. Sin embargo, esta salida es en texto plano sin formato por lo que muchos de los procedimientos recogidos en los menús producen además una salida mucho más comprensible en formato HTML en la ventana de resultados.

Para guardar el contenido de la ventana de resultados en un fichero se utiliza el menú Archivo ▶ Exportar página como HTML. Con esto aparece un cuadro de diálogo en el que hay que indicar el nombre del fichero y la carpeta donde se desea guardar. El fichero resultante está en formato HTML por lo que se podrá visualizar con cualquier navegador web.

7.2 Limpiar la ventana de resultados

La ventana de resultados va acumulando todas las salidas de los análisis realizados en cada sesión de trabajo. Para no mezclar los resultados de estudios distintos, conviene limpiar la ventana de resultados cada vez que se empieza un estudio nuevo. Para ello hay que seleccionar el menú Edición ▶ Limpiar salida.

8 Manipulación de guiones de comandos

8.1 Creación de un guión de comandos

RKWard también incorpora un entorno de desarrollo para programadores de R que permite crear guiones de comandos que pueden ejecutarse todos seguidos. Esta opción es muy interesante para repetir análisis o automatizar tareas repetitivas. Para crear un guión de comandos hay que seleccionar el menú Archivo ▶ Nuevo ▶ Archivo de guiones. Con esto aparecerá una venta como la que aparece en la figura 1.9 donde se podrán teclear los comandos de R para después ejecutarlos uno a uno o en bloque.

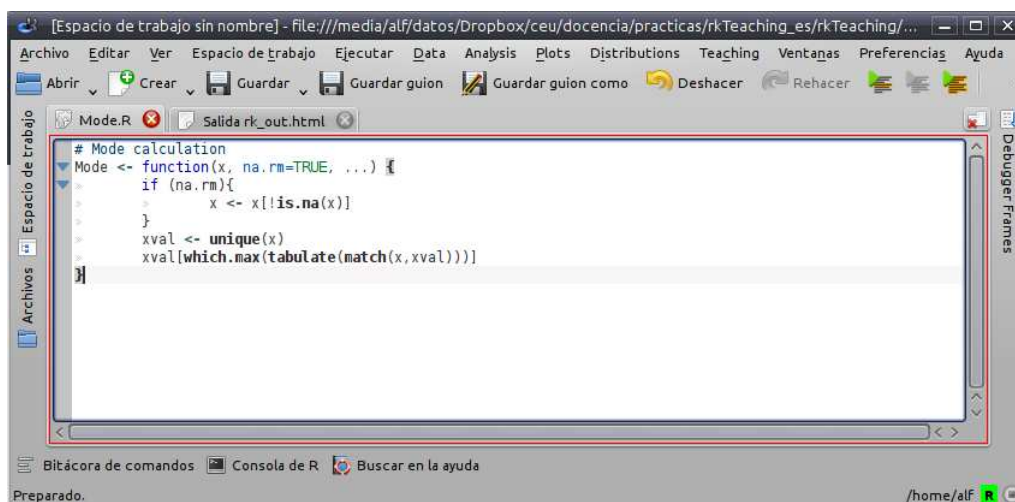


Figura 1.9 – Ventana de edición de guiones de comandos

8.2 Guardar un guión de comandos

Los guiones de comandos también pueden guardarse en un fichero de texto plano mediante el menú Archivo►Guardar guión e indicando el nombre del fichero y la carpeta donde se guardará en el cuadro de diálogo que aparece.

8.3 Abrir un guión de comandos

Para abrir un fichero con un guión de comandos se utiliza el menú Archivo►Abrir archivo de guiones de R y después seleccionar el fichero que se desea abrir en el cuadro de diálogo que aparece.

9 Ayuda

Otra de las ventajas de R es que tiene un sistema de ayuda muy documentado. Es posible conseguir ayuda sobre cualquier función, procedimiento o paquete simplemente tecleando el comando `help()`. Por ejemplo, para obtener ayuda sobre el comando `mean` se teclearía

```
> help("mean")
```

y con esto aparecerá una ventana de ayuda donde se describe la función y también aparecen ejemplos que ilustran su uso. Si no se conoce exactamente el nombre de la función o comando, se puede hacer una búsqueda aproximada con el comando `help.search()`. Por ejemplo, si no se recuerda el nombre de la función logarítmica, se podría teclear

```
> help("logarithm")
```

y con esto aparecerá una ventana con todos los ficheros de ayuda que contienen la palabra `logarithm`.

Finalmente, también es posible invocar la ayuda general de R en RKward con el menú Ayuda►Ayuda de R con lo que aparecerá una página web desde donde podremos navegar a la información deseada. También es posible buscar ayuda sobre un comando concreto en el menú Ayuda►Buscar en la ayuda de R.

Para más información sobre R se recomienda visitar la página <http://www.r-project.org/>, y para más información sobre RKward se recomienda visitar la página <http://rkward.sourceforge.net/>.

10 Ejercicios resueltos

1. Crear un conjunto de datos con los datos de la siguiente muestra y guardarlo con el nombre `colesterol.rda`

Nombre	Sexo	Peso	Altura	Colesterol
José Luis Martínez Izquierdo	H	85	179	182
Rosa Díaz Díaz	M	65	173	232
Javier García Sánchez	H	71	181	191
Carmen López Pinzón	M	65	170	200
Marisa López Collado	M	51	158	148
Antonio Ruiz Cruz	H	66	174	249



Para crear el conjunto de datos:

- a) Seleccionar el menú Archivo►Nuevo►Conjunto de datos.
- b) En el cuadro de diálogo que aparece introducir el nombre del conjunto de datos `colesterol` y hacer clic en el botón Aceptar.
- c) En la ventana del editor de datos hay que definir una variable en cada columna introduciendo su nombre y tipo en las casillas de la cabecera de cada columna.
- d) Una vez definidas las variables hay que introducir los datos de cada variable en la columna correspondiente.

Para guardar los datos:

- a) Seleccionar el menú Espacio de trabajo►Guardar espacio de trabajo.
- b) En el cuadro de diálogo que aparece hay que darle un nombre al fichero, seleccionar la carpeta donde guardarlo y hacer clic en el botón Aceptar.

2. Abrir el fichero creado en el ejercicio anterior y realizar las siguientes operaciones:

- a) Insertar una nueva variable Edad con las edades de todos los individuos de la muestra.

Nombre	Edad
José Luis Martínez Izquierdo	18
Rosa Díaz Díaz	32
Javier García Sánchez	24
Carmen López Pinzón	35
Marisa López Collado	46
Antonio Ruiz Cruz	68



Para abrir el conjunto de datos del ejercicio anterior:

- 1) Seleccionar el menú Espacio de trabajo►Abrir espacio de trabajo.
- 2) En el cuadro de diálogo que aparece seleccionar la carpeta donde se encuentra el fichero con los datos del ejercicio anterior, seleccionar el fichero y hacer clic en el botón Aceptar.

Para insertar la variable Edad:

- 1) Hacer clic en la solapa Espacio de trabajo.
- 2) En la ventana del espacio de trabajo doble clic sobre el conjunto de datos `colesterol`.
- 3) En la ventana del editor de datos introducir el nombre de la variable `edad` y su tipo en las casillas de la cabecera de una nueva columna vacía, e introducir los datos de las edades en las celdas de más abajo.

b) Insertar un nuevo individuo con siguientes datos

Nombre: Cristóbal Campos Ruiz.
 Edad: 44 años.
 Sexo: Hombre.
 Peso: 70 Kg.
 Altura: 178 cm.
 Colesterol: 220 mg/dl.



- 1) En la ventana del editor de datos introducir los datos de del nuevo individuo en la primera fila vacía.

c) Crear una nueva variable donde se calcule el índice de masa corporal de cada paciente mediante la formula:

$$\text{imc} = \frac{\text{Peso (en Kg)}}{\text{Altura (en mt)}^2}$$



- 1) Seleccionar el menú Teaching►Datos►Calcular variable.
- 2) En el cuadro de diálogo que aparece introducir la fórmula para calcular el índice de masa corporal en el campo Expresión de cálculo.
- 3) En el cuadro Guardar nueva variable hacer clic sobre el botón Cambiar.
- 4) En el cuadro de diálogo que aparece seleccionar como objeto padre la el conjunto de datos colesterol y hacer clic sobre el botón Aceptar.
- 5) Introducir el nombre de la nueva variable imc y hacer clic sobre el botón Aceptar.

d) Recodificar el índice de masa corporal en una nueva variable de acuerdo a las siguientes categorías:

Menor de 18,5	Bajo peso
De 18,5 a 24,5	Saludable
De 24,5 a 30	Sobrepeso
Mayor de 30	Obeso



- 1) Seleccionar el menú Teaching►Datos►Recodificar variable.
- 2) En el cuadro de diálogo que aparece seleccionar como variable a recodificar la variable imc.
- 3) Introducir las reglas de recodificación en el campo Reglas de recodificación:

$10:18.5 = 1$
 $18.5:24.5 = 2$
 $24.5:30 = 3$
 $30:hi = 4$
- 4) En el cuadro Guardar nueva variable hacer clic sobre el botón Cambiar.
- 5) En el cuadro de diálogo que aparece seleccionar como objeto padre la el conjunto de datos colesterol y hacer clic sobre el botón Aceptar.
- 6) Introducir el nombre de la nueva variable obesidad y hacer clic sobre el botón Aceptar.
- 7) En la ventana de edición de datos introducir los niveles del factor, asignando Bajo peso a la categoría 1, Saludable a la categoría 2, Sobrepeso a la categoría 3 y Obeso a la categoría 4.

e) Filtrar el conjunto de datos para obtener un nuevo conjunto de datos con los datos de los hombres



- 1) Seleccionar el menú Teaching►Datos►Filtrar.
- 2) En el cuadro de diálogo que aparece seleccionar como conjunto de datos colesterol.
- 3) En el campo Condición de selección introducir la condición `sexo=="H"`.
- 4) Introducir el nombre del nuevo conjunto de datos `colesterol.hombres` y hacer clic sobre el botón Aceptar.

11 Ejercicios propuestos

1. El conjunto de datos neonatos del paquete `rk.Teaching`, contiene información sobre una muestra de 320 recién nacidos en un hospital durante un año que cumplieron el tiempo normal de gestación. Se pide:

- a) Cargar el conjunto de datos.



- 1) Hacer clic en la solapa Espacio de trabajo para desplegarla y ver los paquetes del espacio de trabajo.
- 2) Hacer doble clic sobre el paquete `rk.Teaching` para ver todos los conjuntos de datos que contiene.
- 3) Hacer clic con el botón derecho sobre el conjunto de datos `neonatos` y en el menú contextual que aparece seleccionar Copiar a `.GlobalEnv` para hacer una copia del conjunto de datos en nuestro entorno de trabajo.

- b) Calcular la variable `apgar.medio` como la media de las variables `apgar1` y `apgar5`.
- c) Recodificar la variable `peso` en el factor `categoria.peso` con dos categorías que se correspondan con los pesos menores y mayores de 2,5 Kg.
- d) Recodificar la variable `apgar1` en el factor `estado.apgar1` con tres categorías: deprimido ($\text{Apgar} \leq 3$), moderadamente deprimido ($3 < \text{Apgar} \leq 6$) y normal ($\text{Apgar} > 6$).
- e) Filtrar el conjunto de datos para quedarse con los hijos de las madres no fumadoras con una puntuación Apgar al minuto de nacer menor o igual que 3. ¿Cuántos niños hay?

Distribuciones de Frecuencias y Representaciones Gráficas

1 Fundamentos teóricos

Uno de los primeros pasos en cualquier estudio estadístico es el resumen y la descripción de la información contenida en una muestra. Para ello se van a aplicar algunos métodos de análisis descriptivo, que nos permitirán clasificar y estructurar la información al igual que representarla gráficamente.

Las características que estudiamos pueden ser o no susceptibles de medida; en este sentido definiremos una *variable* como un carácter susceptible de ser medido, es decir, cuantitativo y cuantificable mediante la observación, (por ejemplo el peso de las personas, la edad, etc...), y definiremos un *atributo* como un carácter no susceptible de ser medido, y en consecuencia observable tan sólo cualitativamente (por ejemplo el color de ojos, estado de un paciente, etc...). Se llaman modalidades a las posibles observaciones de un atributo.

Dentro de los atributos, podemos hablar de *atributos ordinales*, los que presentan algún tipo de orden entre las distintas modalidades, y de *atributos nominales*, en los que no existe ningún orden entre ellas.

Dentro de las variables podemos diferenciar entre *discretas*, si sus valores posibles son valores aislados, y *continuas*, si pueden tomar cualquier valor dentro de un intervalo.

En algunos textos no se emplea el término *atributo* y se denominan a todos los caracteres *variables*. En ese caso se distinguen *variables cuantitativas* para designar las que aquí hemos definido como *variables*, y *variables cualitativas* para las que aquí se han llamado *atributos*. En lo sucesivo se aplicará este criterio para simplificar la exposición.

1.1 Cálculo de Frecuencias

Para estudiar cualquier característica, lo primero que deberemos hacer es un recuento de las observaciones, y el número de repeticiones de éstas. Para cada valor x_i de la muestra se define:

Frecuencia absoluta Es el número de veces que aparece cada uno de los valores x_i y se denota por n_i .

Frecuencia relativa Es el número de veces que aparece cada valor x_i dividido entre el tamaño muestral y se denota por f_i

$$f_i = \frac{n_i}{n}$$

Generalmente las frecuencias relativas se multiplican por 100 para que representen el tanto por ciento.

En el caso de que exista un orden entre los valores de la variable, a veces nos interesa no sólo conocer el número de veces que se repite un determinado valor, sino también el número de veces que aparece dicho valor y todos los menores. A este tipo de frecuencias se le denomina *frecuencias acumuladas*.

Frecuencia absoluta acumulada Es la suma de las frecuencias absolutas de los valores menores que x_i más la frecuencia absoluta de x_i , y se denota por N_i

$$N_i = n_1 + n_2 + \dots + n_i$$

Frecuencia relativa acumulada Es la suma de las frecuencias relativas de los valores menores que x_i más la frecuencia relativa de x_i , y se denota por F_i

$$F_i = f_1 + f_2 + \dots + f_i$$

Los resultados de las observaciones de los valores de una variable estadística en una muestra suelen representarse en forma de tabla. En la primera columna se representan los valores x_i de la variable colocados en orden creciente, y en la siguiente columna los valores de las frecuencias absolutas correspondientes n_i .

Podemos completar la tabla con otras columnas, correspondientes a las frecuencias relativas, f_i , y a las frecuencias acumuladas, N_i y F_i . Al conjunto de los valores de la variable observados en la muestra junto con sus frecuencias se le conoce como *distribución de frecuencias muestral*.

■ **Ejemplo 2.1** En una encuesta a 25 matrimonios, sobre el número de hijos que tienen, se obtienen 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.

Los valores distintos de la variable son: 0, 1, 2, 3 y 4. Así la tabla será:

x_i	Recuento	n_i
0	II	2
1	IIII I	6
2	IIII IIII III	14
3	II	2
4	I	1

La distribución de las frecuencias quedaría:

x_i	n_i	f_i	N_i	F_i
0	2	0,08	2	0,08
1	6	0,24	8	0,32
2	14	0,56	22	0,88
3	2	0,08	24	0,96
4	1	0,04	25	1
Suma	25	1		

Cuando el tamaño de la muestra es grande en el caso de variables discretas con muchos valores distintos de la variable, y en cualquier caso si se trata de variables continuas, se agrupan las observaciones en *clases*, que son intervalos contiguos, preferiblemente de la misma amplitud.

Para decidir el número de clases a considerar, una regla frecuentemente utilizada es tomar el entero más próximo a \sqrt{n} donde n es el número de observaciones en la muestra. Pero conviene probar con distintos números de clases y escoger el que proporcione una descripción más clara. Así se prefijan los intervalos $(a_{i-1}, a_i]$, $i = 1, 2, \dots, l$ siendo $a = a_0 < a_1 < \dots < a_l = b$ de tal modo que todos los valores observados estén dentro del intervalo $(a, b]$, y sin que exista ambigüedad a la hora de decidir a qué intervalo pertenece cada dato.

Llamaremos *marca de clase* al punto medio de cada intervalo. Así la *marca de la clase* $(a_{i-1}, a_i]$ es el punto medio x_i de dicha clase, es decir

$$x_i = \frac{a_{i-1} + a_i}{2}$$

En el tratamiento estadístico de los datos agrupados, todos los valores que están en una misma clase se consideran iguales a la marca de la clase. De esta manera si en la clase $(a_{i-1}, a_i]$ hay n_i valores observados, se puede asociar la marca de la clase x_i con esta frecuencia n_i .

1.2 Representaciones Gráficas

Hemos visto que la tabla estadística resume los datos de una muestra, de forma que ésta se puede analizar de una manera más sistemática y resumida. Para conseguir una percepción visual de las características de la población resulta muy útil el uso de gráficas y diagramas. Dependiendo del tipo de variable y de si trabajamos con datos agrupados o no, se utilizarán distintos tipos.

Diagrama de barras y polígono de frecuencias

Consiste en representar sobre el eje de abscisas de un sistema de ejes coordenados los distintos valores de la variable X , y levantar sobre cada uno de esos puntos una barra cuya altura sea igual a la frecuencia absoluta o relativa correspondiente a ese valor, tal y como se muestra en la figura 2.1(a). Esta representación se utiliza para distribuciones de frecuencias con pocos valores distintos de la variable, tanto cuantitativas como cualitativas, y en este último caso se suele representar con rectángulos de altura igual a la frecuencia de cada modalidad.

En el caso de variables cuantitativas se puede representar también el diagrama de barras de las frecuencias acumuladas, tal y como se muestra en la figura 2.1(b).

Otra representación habitual es el *polígono de frecuencias* que consiste en la línea poligonal cuyos vertices son los puntos (x_i, n_i) , tal y como se ve en la figura 2.1(c), y si en vez de considerar las frecuencias absolutas o relativas se consideran las absolutas o relativas acumuladas, se obtiene el *polígono de frecuencias acumuladas*, como se ve en la figura 2.1(d).

Histogramas

Este tipo de representaciones se utiliza en variables continuas y en variables discretas en que se ha realizado una agrupación de las observaciones en clases. Un *histograma* es un conjunto de rectángulos, cuyas bases son los intervalos de clase $(a_{i-1}, a_i]$ sobre el eje OX y su altura la correspondiente frecuencia absoluta, relativa, absoluta acumulada, o relativa acumulada, tal y como se muestra en la figuras 2.2(a) y 2.2(b).

Si unimos los puntos medios de las bases superiores de los rectángulos del histograma, se obtiene el *polígono de frecuencias* correspondiente a datos agrupados (figura 2.2(c)).

El polígono de frecuencias también se puede utilizar para representar las frecuencias acumuladas, tanto absolutas como relativas. En este caso la línea poligonal se traza uniendo los extremos derechos de las bases superiores de los rectángulos del histograma de frecuencias acumuladas, en lugar de los puntos centrales (figura 2.2(d)).

Para variables cualitativas y cuantitativas discretas también se pueden usar las superficies representativas; de éstas, las más empleadas son los *sectores circulares*.

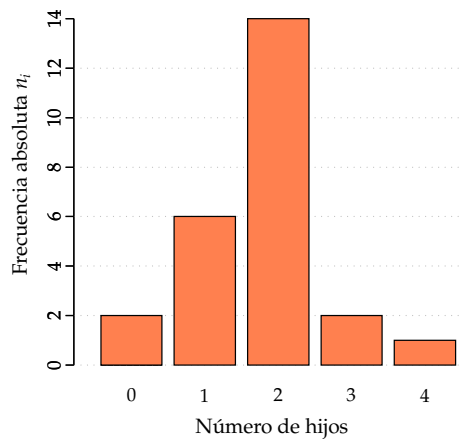
Sectores circulares o diagrama de sectores

Es una representación en la que un círculo se divide en sectores, de forma que los ángulos, y por tanto las áreas respectivas, sean proporcionales a la frecuencia.

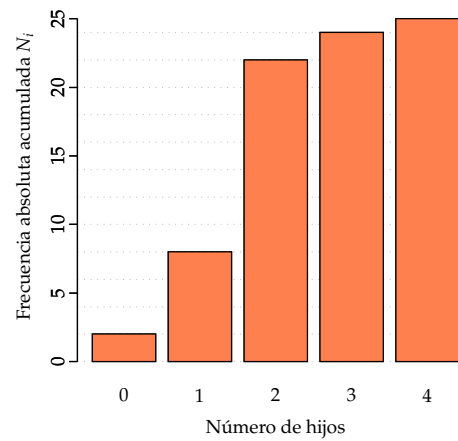
■ **Ejemplo 2.2** Se está haciendo un estudio en una población del grupo sanguíneo de sus ciudadanos. Para ello disponemos de una muestra de 30 personas, con los siguientes resultados: 5 personas con grupo 0, 14 con grupo A, 8 con grupo B y 3 con grupo AB. El diagrama de sectores de frecuencias relativas correspondiente aparece en la figura 2.3.

Diagrama de cajas y datos atípicos

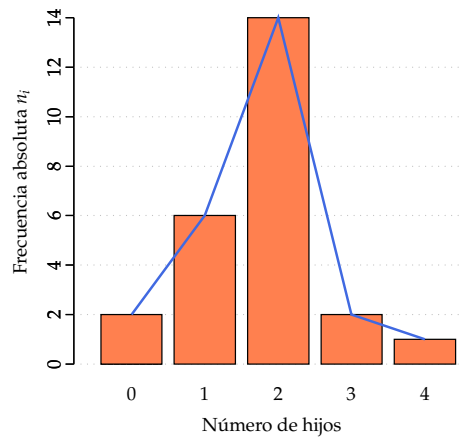
Los datos extremadamente altos o bajos, en comparación con los del resto de la muestra, reciben el nombre de datos influyentes o *datos atípicos*. Tales datos que, como su propio nombre indica, pueden modificar las conclusiones de un estudio, deben ser considerados atentamente antes de aceptarlos, pues no pocas veces podrán ser, simplemente, datos erróneos. La representación gráfica más apropiada para detectar estos datos es el *diagrama de cajas*. Este diagrama está formado por una caja que contiene el 50 % de los datos centrales de la distribución, y unos segmentos que salen de la caja, que



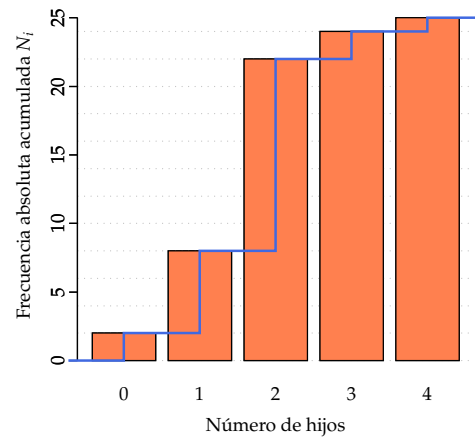
(a) Diagrama de barras de frecuencias absolutas.



(b) Diagrama de barras de frecuencias absolutas acumuladas.



(c) Polígono de frecuencias absolutas.

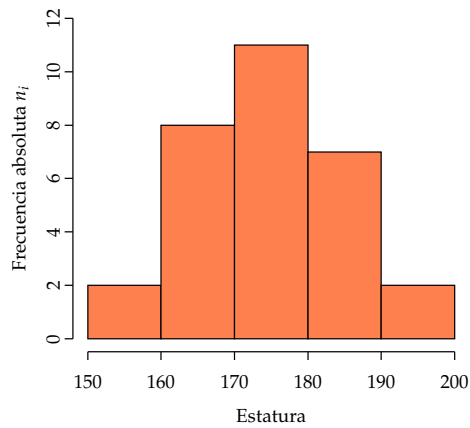


(d) Polígono de frecuencias absolutas acumuladas

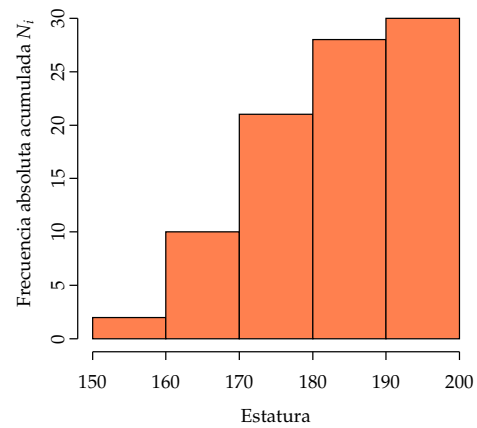
Figura 2.1 – Diagramas de barras y polígonos asociados para datos no agrupados.

indican los límites a partir de los cuales los datos se consideran atípicos. En la figura 2.4 se puede observar un ejemplo en el que aparecen dos datos atípicos.

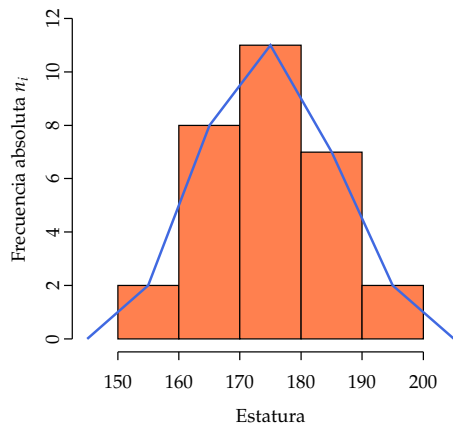
2. Distribuciones de Frecuencias y Representaciones Gráficas



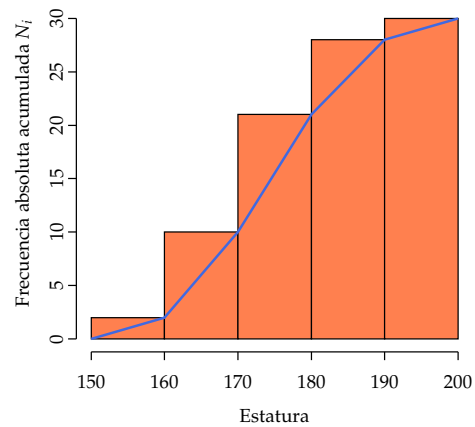
(a) Histograma de frecuencias absolutas.



(b) Histograma de frecuencias absolutas acumuladas.



(c) Polígono de frecuencias absolutas.



(d) Polígono de frecuencias absolutas acumuladas.

Figura 2.2 – Histograma y polígonos asociados para datos agrupados.

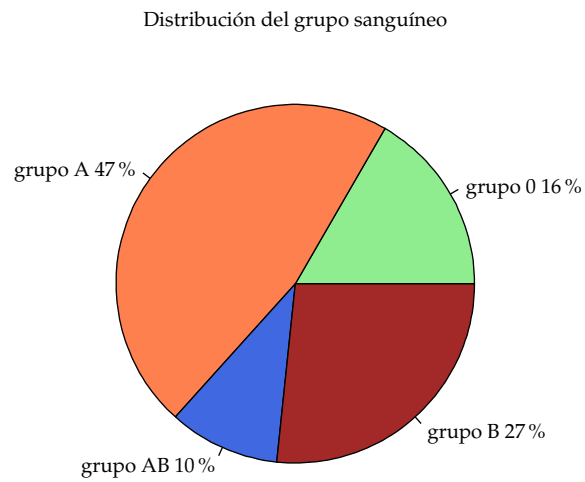


Figura 2.3 – Diagrama de sectores de frecuencias relativas del grupo sanguíneo.

Diagrama de caja y bigotes del peso de recién nacidos

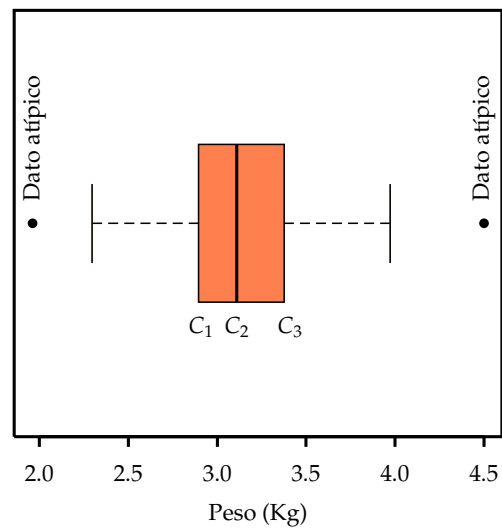


Figura 2.4 – Diagrama de cajas para una muestra de recién nacidos. Existen dos niños con pesos atípicos, uno con peso extremadamente bajo 1,9 kg, y otro con peso extremadamente alto 4,3 kg.

2 Ejercicios resueltos

1. En una encuesta a 25 matrimonios sobre el número de hijos que tenían se obtuvieron 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

Se pide:

- a) Crear un conjunto de datos con la variable hijos e introducir los datos.
- b) Construir la tabla de frecuencias.



- 1) Seleccionar el menú Teaching►Distribución de frecuencias►Tabla de frecuencias.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable hijos en el campo Variable a tabular y hacer clic en el botón Enviar.

- c) Dibujar el diagrama de barras de las frecuencias absolutas.



- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de barras.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable hijos en el campo Variable y hacer clic en el botón Enviar.

- d) Para la misma tabla de frecuencias anterior, dibujar también el diagrama de barras de las frecuencias relativas, el de absolutas acumuladas y el de relativas acumuladas, además de sus correspondientes polígonos.



Repetir los pasos del apartado anterior activando, en la solapa de Opciones de las barras, la opción Frecuencias relativas si se desea el diagrama de barras de frecuencias relativas, activando la opción Frecuencias acumuladas si se desea el diagrama de barras de frecuencias acumuladas y activando la opción Polígono para obtener el polígono asociado.

2. En un hospital se realizó un estudio sobre el número de personas que ingresaron en urgencias cada día del mes de noviembre. Los datos observados fueron:

15, 23, 12, 10, 28, 50, 12, 17, 20, 21, 18, 13, 11, 12, 26
30, 6, 16, 19, 22, 14, 17, 21, 28, 9, 16, 13, 11, 16, 20

Se pide:

- a) Crear un conjunto de datos con la variable urgencias e introducir los datos.
- b) Dibujar el diagrama de cajas. ¿Existe algún dato atípico? En el caso de que exista, eliminarlo y proceder con los siguientes apartados.



- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de cajas.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable urgencias en el campo Variables y hacer clic en el botón Enviar.
- 3) En la ventana que aparece con el diagrama de cajas identificar el dato atípico.
- 4) Ir a la ventana de edición de datos y eliminar la fila del dato atípico haciendo clic con el botón derecho del ratón en la cabecera de la fila y seleccionando Borrar esta fila.

- c) Construir la tabla de frecuencias agrupando en 5 clases.



- 1) Seleccionar el menú Teaching►Distribución de frecuencias►Tabla de frecuencias.
- 2) En el cuadro de diálogo que aparece seleccionar la variable urgencias.
- 3) En la solapa de Clases activar la casilla Agrupar en intervalos, marcar la opción Número de intervalos e introducir el número deseado de intervalos en el campo Intervalos sugeridos y hacer clic sobre el botón Enviar.

d) Dibujar el histograma de frecuencias absolutas correspondiente a la tabla anterior.



- 1) Seleccionar el menú Teaching►Gráficos►Histograma.
- 2) En el cuadro de diálogo que aparece seleccionar la variable urgencias en el campo Variable.
- 3) En la solapa de Clases activar la casilla Agrupar en intervalos, marcar la opción Número de intervalos e introducir el número deseado de intervalos en el campo Intervalos sugeridos y hacer clic sobre el botón Enviar.

e) Para la misma tabla de frecuencias anterior, dibujar también el histograma de las frecuencias relativas, el de absolutas acumuladas y el de relativas acumuladas, además de sus correspondientes polígonos.



Repetir los pasos del apartado anterior activando, en la solapa de Opciones del histograma, la opción Frecuencias relativas si se desea el histograma de frecuencias relativas, activando la opción Frecuencias acumuladas si se desea el histograma de frecuencias acumuladas y activando la opción Polígono para obtener el polígono asociado.

3. Los grupos sanguíneos de una muestra de 30 personas son:

A, B, B, A, AB, 0, 0, A, B, B, A, A, A, A, AB,
A, A, A, B, 0, B, B, B, A, A, A, 0, A, AB, 0.

Se pide:

- a) Crear un conjunto de datos con la variable grupo.sanguineo e introducir los datos.
- b) Construir la tabla de frecuencias.



- 1) Seleccionar el menú Teaching►Distribución de frecuencias►Tabla de frecuencias .
- 2) En el cuadro de diálogo que aparece, seleccionar la variable grupo.sanguineo en el campo Variable a tabular y hacer clic en el botón Enviar.

c) Dibujar el diagrama de sectores.



- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de sectores.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable grupo.sanguineo en el campo Variables y hacer clic sobre el botón Enviar.

4. En un estudio de población se tomó una muestra de 27 personas, y se les preguntó por su edad y estado civil, obteniendo los siguientes resultados:

Estado civil	Edad									
Soltero	31	45	35	65	21	38	62	22	31	
Casado	62	39	62	59	21	62				
Viudo	80	68	65	40	78	69	75			
Divorciado	31	65	59	49	65					

Se pide:

- Crear un conjunto de datos con la variables estado.civil y edad e introducir los datos.
- Construir la tabla de frecuencias de la variable edad para cada categoría de la variable estado.civil.



- Seleccionar el menú Teaching►Distribución de frecuencias►Tabla de frecuencias.
- En el cuadro de diálogo que aparece, seleccionar la variable edad en el campo Variable a tabular, activar la casilla Tabular por grupos, seleccionar la variable estado.civil en el campo Variable de agrupación y hacer clic en el botón Enviar.

- Dibujar los diagramas de cajas de la edad según el estado civil. ¿Existen datos atípicos? ¿En qué grupo hay mayor dispersión?



- Seleccionar el menú Teaching►Gráficos►Diagrama de cajas.
- En el cuadro de diálogo que aparece, seleccionar la variable edad en el campo Variables, activar la casilla Dibujar por grupos, seleccionar la variable estado.civil en el campo Variable de agrupación y hacer clic en el botón Enviar.

3 Ejercicios propuestos

- El número de lesiones padecidas durante una temporada por cada jugador de un equipo de fútbol fue el siguiente:

0, 1, 2, 1, 3, 0, 1, 0, 1, 2, 0, 1, 1, 1, 2, 0, 1, 3, 2, 1, 2, 1, 0, 1

Se pide:

- Construir la tabla de frecuencias.
 - Dibujar el diagrama de barras de las frecuencias relativas y de frecuencias relativas acumuladas.
 - Dibujar el diagrama de sectores.
- Para realizar un estudio sobre la estatura de los estudiantes universitarios, seleccionamos, mediante un proceso de muestreo aleatorio, una muestra de 30 estudiantes, obteniendo los siguientes resultados (medidos en centímetros):

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.

Se pide:

- Dibujar el histograma de las frecuencias absolutas agrupando desde 150 a 200 en clases de amplitud 10.
 - Dibujar el diagrama de cajas. ¿Existe algún dato atípico?.
- El conjunto de datos neonatos del paquete rk.Teaching, contiene información sobre una muestra de 320 recién nacidos en un hospital durante un año que cumplieron el tiempo normal de gestación. Se pide:
 - Construir la tabla de frecuencias de la puntuación Apgar al minuto de nacer. Si se considera que una puntuación Apgar de 3 o menos indica que el neonato está deprimido, ¿qué porcentaje de niños está deprimido en la muestra?

- b) Comparar las distribuciones de frecuencias de las puntuaciones Apgar al minuto de nacer según si la madre es mayor o menor de 20 años. ¿En qué grupo hay más neonatos deprimidos?
- c) Construir la tabla de frecuencias para el peso de los neonatos, agrupando en clases de amplitud 0,5 desde el 2 hasta el 4,5. ¿En qué intervalo de peso hay más niños?
- d) Comparar la distribución de frecuencias relativas del peso de los neonatos según si la madre fuma o no. Si se considera como peso bajo un peso menor de 2,5 kg, ¿En qué grupo hay un mayor porcentaje de niños con peso bajo?
- e) Si en los recién nacidos se considera como peso bajo un peso menor de 2,5 kg, calcular la prevalencia del bajo peso de recién nacidos en el grupo de madres fumadoras y en el de no fumadoras.
- f) Calcular el riesgo relativo de que un recién nacido tenga bajo peso cuando la madre fuma, frente a cuando la madre no fuma.
- g) Construir el diagrama de barras de la puntuación Apgar al minuto. ¿Qué puntuación Apgar es la más frecuente?
- h) Construir el diagrama de frecuencias relativas acumuladas de la puntuación Apgar al minuto. ¿Por debajo de qué puntuación estarán la mitad de los niños?
- i) Comparar mediante diagramas de barras de frecuencias relativas las distribuciones de las puntuaciones Apgar al minuto según si la madre ha fumado o no durante el embarazo. ¿Qué se puede concluir?
- j) Construir el histograma de pesos, agrupando en clases de amplitud 0,5 desde el 2 hasta el 4,5. ¿En qué intervalo de peso hay más niños?
- k) Comparar la distribución de frecuencias relativas del peso de los neonatos según si la madre fuma o no. ¿En qué grupo se aprecia menor peso de los niños de la muestra?
- l) Comparar la distribución de frecuencias relativas del peso de los neonatos según si la madre fumaba o no antes del embarazo. ¿Qué se puede concluir?
- m) Construir el diagrama de caja y bigotes del peso. ¿Entre qué valores se considera que el peso de un neonato es normal? ¿Existen datos atípicos?
- n) Comparar el diagrama de cajas y bigotes del peso, según si la madre fumó o no durante el embarazo y si era mayor o no de 20 años. ¿En qué grupo el peso tiene más dispersión central? ¿En qué grupo pesan menos los niños de la muestra?
- ñ) Comparar el diagrama de cajas de la puntuación Apgar al minuto y a los cinco minutos. ¿En qué variable hay más dispersión central?

Estadísticos Muestrales

1 Fundamentos teóricos

Hemos visto cómo podemos presentar la información que obtenemos de la muestra, a través de tablas o bien a través de gráficas. La tabla de frecuencias contiene toda la información de la muestra pero resulta difícil sacar conclusiones sobre determinados aspectos de la distribución con sólo mirarla. Ahora veremos cómo a partir de esos mismos valores observados de la variable estadística, se calculan ciertos números que resumen la información muestral. Estos números, llamados *Estadísticos*, se utilizan para poner de manifiesto ciertos aspectos de la distribución, tales como la dispersión o concentración de los datos, la forma de su distribución, etc. Según sea la característica que pretenden reflejar se pueden clasificar en medidas de posición, medidas de dispersión y medidas de forma.

1.1 Medidas de posición

Son valores que indican cómo se sitúan los datos. Los más importantes son la Media aritmética, la Mediana y la Moda.

Media aritmética \bar{x}

Se llama *media aritmética* de una variable estadística X , y se representa por \bar{x} , a la suma de todos los resultados observados, dividida por el tamaño muestral. Es decir, la media de la variable estadística X , cuya distribución de frecuencias es (x_i, n_i) , viene dada por

$$\bar{x} = \frac{x_1 + \dots + x_1 + \dots + x_k + \dots + x_k}{n_1 + \dots + n_k} = \frac{x_1 n_1 + \dots + x_k n_k}{n} = \frac{1}{n} \sum_{i=1}^k x_i n_i$$

La media aritmética sólo tiene sentido en variables cuantitativas.

Mediana Me

Se llama *mediana* y lo denotamos por Me , a aquel valor de la muestra que, una vez ordenados todos los valores de la misma en orden creciente, tiene tantos términos inferiores a él como superiores. En consecuencia, divide la distribución en dos partes iguales.

La mediana sólo tiene sentido en atributos ordinales y en variables cuantitativas.

Moda Mo

La *moda* es el valor de la variable que presenta una mayor frecuencia en la muestra. Cuando haya más de un valor con frecuencia máxima diremos que hay más de una moda. En variables continuas o discretas agrupadas llamaremos clase modal a la que tenga la máxima frecuencia. Se puede calcular la moda tanto en variables cuantitativas como cualitativas.

Cuantiles

Si el conjunto total de valores observados se divide en r partes que contengan cada una $\frac{n}{r}$ observaciones, los puntos de separación de las mismas reciben el nombre genérico de *cuantiles*.

Según esto la mediana también es un cuantil con $r = 2$. Algunos cuantiles reciben determinados nombres como:

Cuartiles. Son los puntos que dividen la distribución en 4 partes iguales y se designan por C_1, C_2, C_3 . Es claro que $C_2 = Me$.

Deciles. Son los puntos que dividen la distribución en 10 partes iguales y se designan por D_1, D_2, \dots, D_9 .

Percentiles. Son los puntos que dividen la distribución en 100 partes iguales y se designan por P_1, P_2, \dots, P_{99} .

1.2 Medidas de dispersión

Miden la separación existente entre los valores de la muestra. Las más importantes son el Rango o Recorrido, el Rango Intercuartílico, la Varianza, la Desviación Típica y el Coeficiente de Variación.

Rango o Recorrido Re

La medida de dispersión más inmediata es el rango. Llamamos *recorrido* o *rango* y lo designaremos por Re a la diferencia entre los valores máximo y mínimo que toma la variable en la muestra, es decir

$$Re = \max\{x_i, i = 1, 2, \dots, n\} - \min\{x_i, i = 1, 2, \dots, n\}.$$

Este estadístico sirve para medir el campo de variación de la variable, aunque es la medida de dispersión que menos información proporciona sobre la mayor o menor agrupación de los valores de la variable alrededor de las medidas de tendencia central. Además tiene el inconveniente de que se ve muy afectado por los datos atípicos.

Rango Intercuartílico RI

El *rango intercuartílico* RI es la diferencia entre el tercer y el primer cuartil, y mide, por tanto, el campo de variación del 50 % de los datos centrales de la distribución. Por consiguiente

$$RI = C_3 - C_1.$$

La ventaja del rango intercuartílico frente al recorrido es que no se ve tan afectado por los datos atípicos.

Varianza s_x^2

Llamamos *varianza* de una variable estadística X , y la designaremos por s_x^2 , a la media de los cuadrados de las desviaciones de los valores observados respecto de la media de la muestra, es decir,

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

Desviación Típica s_x

La raíz cuadrada positiva de la varianza se conoce como *desviación típica* de la variable X , y se representa por s ,

$$s = +\sqrt{s_x^2}.$$

Coeficiente de Variación de Pearson Cv_x

Al cociente entre la desviación típica y el valor absoluto de la media se le conoce como *coeficiente de variación de Pearson* o simplemente *coeficiente de variación*:

$$Cv_x = \frac{s_x}{|\bar{x}|}.$$

El coeficiente de variación es adimensional, y por tanto permite hacer comparaciones entre variables expresadas en distintas unidades. Cuanto más próximo esté a 0, menor será la dispersión de la muestra en relación con la media, y más representativa será ésta última del conjunto de observaciones.

1.3 Medidas de forma

Indican la forma que tiene la distribución de valores en la muestra. Se pueden clasificar en dos grupos: Medidas de *asimetría* y medidas de *apuntamiento o curtosis*.

Coeficiente de asimetría de Fisher g_1

El *coeficiente de asimetría de Fisher*, que se representa por g_1 , se define

$$g_1 = \frac{\sum_{i=1}^k (x_i - \bar{x})^3 f_i}{s_x^3}.$$

Dependiendo del valor que tome tendremos:

- $g_1 = 0$. Distribución simétrica.
- $g_1 < 0$. Distribución asimétrica hacia la izquierda.
- $g_1 > 0$. Distribución asimétrica hacia la derecha.

Coeficiente de apuntamiento o curtosis g_2

El grado de apuntamiento de las observaciones de la muestra, se caracteriza por el *coeficiente de apuntamiento o curtosis*, que se representa por g_2 , y se define

$$g_2 = \frac{\sum_{i=1}^k (x_i - \bar{x})^4 f_i}{s_x^4} - 3.$$

Dependiendo del valor que tome tendremos:

- $g_2 = 0$. La distribución tiene un apuntamiento igual que el de la distribución normal de la misma media y desviación típica. Se dice que es una distribución *mesocúrtica*.
- $g_2 < 0$. La distribución es menos apuntada que la distribución normal de la misma media y desviación típica. Se dice que es una distribución *platicúrtica*.
- $g_2 > 0$. La distribución es más apuntada que la distribución normal de la misma media y desviación típica. Se dice que es una distribución *leptocúrtica*.

Tanto g_1 como g_2 suelen utilizarse para comprobar si los datos muestrales provienen de una población no normal. Cuando g_1 está fuera del intervalo $[-2,2]$ se dice que la distribución es demasiado asimétrica como para que los datos provengan de una población normal. Del mismo modo, cuando g_2 está fuera del intervalo $[-2,2]$ se dice que la distribución es, o demasiado apuntada, o demasiado plana, como para que los datos provengan de una población normal.

1.4 Estadísticos de variables en las que se definen grupos

Ya sabemos cómo resumir la información contenida en una muestra utilizando una serie de estadísticos. Pero hasta ahora sólo hemos estudiado ejemplos con un único carácter objeto de estudio.

En la mayoría de las investigaciones no estudiaremos un único carácter, sino un conjunto de caracteres, y muchas veces será conveniente obtener información de un determinado carácter, en función de los grupos creados por otro de los caracteres estudiados en la investigación. A estas variables que se utilizan para formar grupos se les conoce como *variables clasificadoras* o *factores*.

Por ejemplo, si se realiza un estudio sobre un conjunto de niños recién nacidos, podemos estudiar su peso. Pero si además sabemos si la madre de cada niño es fumadora o no, podremos hacer un estudio del peso de los niños de las madres fumadoras por un lado y los de las no fumadoras por otro, para ver si existen diferencias entre ambos grupos.

2 Ejercicios resueltos

1. En una encuesta a 25 matrimonios sobre el número de hijos que tenían se obtuvieron 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

Se pide:

- Crear un conjunto de datos con la variable hijos e introducir los datos. Si ya se tienen los datos, simplemente recuperarlos.
- Calcular la media aritmética, varianza y desviación típica de dicha variable. Interpretar los estadísticos.



- Seleccionar el menú Teaching►Estadística descriptiva►Estadísticos.
- En el cuadro de diálogo que aparece seleccionar la variable hijos en el campo Variables.
- En la solapa Estadísticos básicos seleccionar Media y Desviación típica, y hacer click sobre el botón Enviar.

- Calcular los cuartiles, el recorrido, el rango intercuartílico, el tercer decil y el percentil 68.



- Seleccionar el menú Teaching►Estadística descriptiva►Estadísticos.
- En el cuadro de diálogo que aparece seleccionar la variable hijos en el campo Variables.
- En la solapa Estadísticos básicos seleccionar Cuartiles, Rango, Rango intercuartílico, introducir los valores 0,3 y 0,68 en el campo Percentiles, y hacer click sobre el botón Enviar.

2. En un hospital se realizó un estudio sobre el número de personas que ingresaron en urgencias cada día del mes de noviembre. Los datos observados fueron:

15, 23, 12, 10, 28, 50, 12, 17, 20, 21, 18, 13, 11, 12, 26
30, 6, 16, 19, 22, 14, 17, 21, 28, 9, 16, 13, 11, 16, 20

Se pide:

- Crear un conjunto de datos con la variable urgencias e introducir los datos.
- Calcular la media aritmética, varianza, desviación típica y coeficiente de variación de dicha variable. Interpretar los estadísticos.



- Seleccionar el menú Teaching►Estadística descriptiva►Estadísticos.
- En el cuadro de diálogo que aparece seleccionar la variable urgencias en el campo Variables.
- En la solapa Estadísticos básicos seleccionar Media, Varianza, Desviación típica y Coeficiente de variación, y hacer click sobre el botón Enviar.

- Calcular el coeficiente de asimetría y el de curtosis e interpretar los resultados



Seguir los mismos pasos del apartado anterior, seleccionando Coeficiente de asimetría y Coeficiente de Curtosis en la solapa Estadísticos básicos.

3. En un grupo de 20 alumnos, las calificaciones obtenidas en Matemáticas fueron:

SS, AP, SS, AP, AP, NT, NT, AP, SB, SS
SB, SS, AP, AP, NT, AP, SS, NT, SS, NT

Se pide:

- a) Crear un conjunto de datos curso con la variable calificaciones e introducir los datos.
- b) Recodificar esta variable, asignando 2,5 al SS, 6 al AP, 8 al NT y 9,5 al SB.



- 1) Seleccionar el menú Teaching►Datos►Recodificar variable.
- 2) En el cuadro de diálogo que aparece seleccionar como variable a recodificar la variable calificaciones.
- 3) Introducir las reglas de recodificación en el campo Reglas de recodificación:


```
"SS" = 2.5
"AP" = 6
"NT" = 8
"SB" = 9.5
```
- 4) En el cuadro Guardar nueva variable hacer click sobre el botón Cambiar.
- 5) En el cuadro de diálogo que aparece seleccionar como objeto padre la el conjunto de datos curso y hacer click sobre el botón Enviar.
- 6) Introducir el nombre de la nueva variable nota, desmarcar la casilla Convertir en factor y hacer click sobre el botón Enviar.

- c) La mediana y el rango intercuartílico.



- 1) Seleccionar el menú Teaching►Estadística descriptiva►Estadísticos.
- 2) En el cuadro de diálogo que aparece seleccionar la variable nota en el campo Variables.
- 3) En la solapa Estadísticos básicos seleccionar Mediana y Rango intercuartílico, y hacer click sobre el botón Enviar.

4. Para realizar un estudio sobre la estatura de los estudiantes universitarios se ha seleccionado mediante un proceso de muestreo aleatorio, una muestra de 30 estudiantes, obteniendo los siguientes resultados (medidos en centímetros):

Mujeres: 173, 158, 174, 166, 162, 177, 165, 154, 166, 182, 169, 172, 170, 168.

Hombres: 179, 181, 172, 194, 185, 187, 198, 178, 188, 171, 175, 167, 186, 172, 176, 187.

Se pide:

- a) Crear un conjunto de datos con las variables estatura y sexo e introducir los datos.
- b) Obtener un resumen de estadísticos en el que se muestren la media aritmética, mediana, varianza, desviación típica y cuartiles según el sexo. Interpretar los estadísticos.



- 1) Seleccionar el menú Teaching►Estadística descriptiva►Estadísticos.
- 2) En el cuadro de diálogo que aparece seleccionar la variable estatura en el campo Variables, marcar la casilla Estadística por grupos y seleccionar la variable sexo en el campo Variables de agrupación.
- 3) En la solapa Estadísticos básicos seleccionar Media, Mediana, Varianza, Desviación típica y Cuartiles, y hacer click sobre el botón Enviar.

3 Ejercicios propuestos

1. El número de lesiones padecidas durante una temporada por cada jugador de un equipo de fútbol fue el siguiente:

0, 1, 2, 1, 3, 0, 1, 0, 1, 2, 0, 1, 1, 1, 2, 0, 1, 3, 2, 1, 2, 1, 0, 1

Se pide:

- Calcular la media aritmética, mediana, varianza y desviación típica de las lesiones e interpretarlas.
 - Calcular los coeficientes de asimetría y curtosis e interpretarlos.
 - Calcular el cuarto y el octavo decil e interpretarlos.
2. En un estudio de población se tomó una muestra de 27 personas, y se les preguntó por su edad y estado civil, obteniendo los siguientes resultados:

Estado civil	Edad									
Soltero	31	45	35	65	21	38	62	22	31	
Casado	62	39	62	59	21	62				
Viudo	80	68	65	40	78	69	75			
Divorciado	31	65	59	49	65					

Se pide:

- Calcular la media y la desviación típica de la edad según el estado civil e interpretarlas.
 - ¿En qué grupo es más representativa la media?
3. En un estudio se ha medido la tensión arterial de 25 individuos. Además se les ha preguntado si fuman y beben:

Fumador	si	no	si	si	si	no	no	si	no	si	no	si	no
Bebedor	no	no	si	si	no	no	si	si	no	si	no	si	si
Tensión arterial	80	92	75	56	89	93	101	67	89	63	98	58	91

Fumador	si	no	no	si	no	no	no	si	no	si	no	si	
Bebedor	si	no	si	si	no	no	si	si	si	no	si	no	
Tensión arterial	71	52	98	104	57	89	70	93	69	82	70	49	

Calcular la media aritmética, desviación típica, coeficiente de asimetría y curtosis de la tensión arterial por grupos dependiendo de si beben o fuman e interpretarlos.

4. El conjunto de datos neonatos del paquete `rk.Teaching`, contiene información sobre una muestra de 320 recién nacidos en un hospital durante un año que cumplieron el tiempo normal de gestación. Se pide:
- Calcular la media y la mediana muestral del peso de los nacidos e interpretarlos.
 - Calcular el peso medio de los recién nacidos de la muestra según si la madre ha fumado o no durante el embarazo. Calcular también el peso medio de los recién nacidos de madres que no han fumado durante el embarazo, según si la madre fumaba o no antes del embarazo. ¿Qué conclusiones se pueden sacar?
 - ¿Cuál es la puntuación Apgar al minuto de nacer más frecuente?
 - Calcular la media de la diferencia entre las puntuaciones Apgar a los 5 minutos y al minuto de nacer. ¿Cómo evolucionan los recién nacidos?
 - Calcular los cuartiles muestrales del peso de los recién nacidos e interpretarlos.
 - Comparar los cuartiles muestrales del peso de los recién nacidos según el sexo.
 - ¿Por encima de qué peso estarán el 10 % de los niños con mayor peso?
 - Si se considera que un niño es atípico por bajo peso si se encuentra entre el 5 % de los pesos más bajos, ¿por debajo de qué peso tiene que estar?
 - Calcular el recorrido y el rango intercuartílico muestrales del peso de los recién nacidos e interpretarlos.

- j) Calcular la varianza y la desviación típica del peso de los recién nacidos e interpretarlos.
 - k) ¿En qué grupo hay más variabilidad del peso de los recién nacidos, en las madres fumadoras o en las madres no fumadoras durante el embarazo? ¿En qué grupo será más representativo el peso medio?
 - l) ¿Qué variable presenta más variabilidad relativa, el peso de los recién nacidos o el Apgar al minuto de nacer?
 - m) Calcular el coeficiente de asimetría y de apuntamiento muestrales del peso de los recién nacidos e interpretarlos.
 - n) ¿Qué distribución es más asimétrica, la de los pesos de recién nacidos en madres mayores de 20 años o en madres menores de 20 años?
 - ñ) ¿Qué distribución es más apuntada, la del peso de los recién nacidos en hombres o en mujeres?
 - o) De acuerdo a la forma de la distribución, ¿puede considerarse la puntuación Apgar al minuto de nacer como una variable normal? ¿Y el número de cigarros fumados al día durante el embarazo?
5. Se quiere comparar la precisión de dos tensiómetros, uno de brazo y otro de muñeca, y para ello se han realizado 8 medidas repetidas de la tensión arterial de una misma persona con cada uno de ellos, obteniendo los siguientes valores en mmHg:
- tens.brazo: 111, 109, 112, 111, 113, 113, 114, 111.
 - tens.muñeca: 115, 113, 117, 116, 112, 112, 117, 112.
- ¿Qué tensiómetro es más preciso?

Regresión Lineal Simple y Correlación

1 Fundamentos teóricos

1.1 Regresión

La *regresión* es la parte de la estadística que trata de determinar la posible relación entre una variable numérica Y , que suele llamarse *variable dependiente*, y otro conjunto de variables numéricas, X_1, X_2, \dots, X_n , conocidas como *variables independientes*, de una misma población. Dicha relación se refleja mediante un modelo funcional $y = f(x_1, \dots, x_n)$.

El caso más sencillo se da cuando sólo hay una variable independiente X , y entonces se habla de *regresión simple*. En este caso el modelo que explica la relación entre X e Y es una función de una variable $y = f(x)$.

Dependiendo de la forma de esta función, existen muchos tipos de regresión simple. Los más habituales son los que aparecen en la siguiente tabla:

Modelo	Ecuación genérica
Lineal	$y = a + bx$
Parabólico	$y = a + bx + cx^2$
Polinómico de grado n	$y = a_0 + a_1x + \dots + a_nx^n$
Potencial	$y = ax^b$
Exponencial	$y = e^{a+bx}$
Logarítmico	$y = a + b \log x$
Inverso	$y = a + b/x$
Curva S	$y = e^{a+b/x}$

Para elegir un tipo de modelo u otro, se suele representar el *diagrama de dispersión*, que consiste en dibujar sobre unos ejes cartesianos correspondientes a las variables X e Y , los pares de valores (x_i, y_j) observados en cada individuo de la muestra.

■ **Ejemplo 4.1** En la figura la figura 4.1 aparece el diagrama de dispersión correspondiente a una muestra de 30 individuos en los que se ha medido la estatura en cm (X) y el peso en kg (Y). En este caso la forma de la nube de puntos refleja una relación lineal entre la estatura y el peso.

Según la forma de la nube de puntos del diagrama, se elige el modelo más apropiado (figura 4.2), y se determinan los parámetros de dicho modelo para que la función resultante se ajuste lo mejor posible a la nube de puntos.

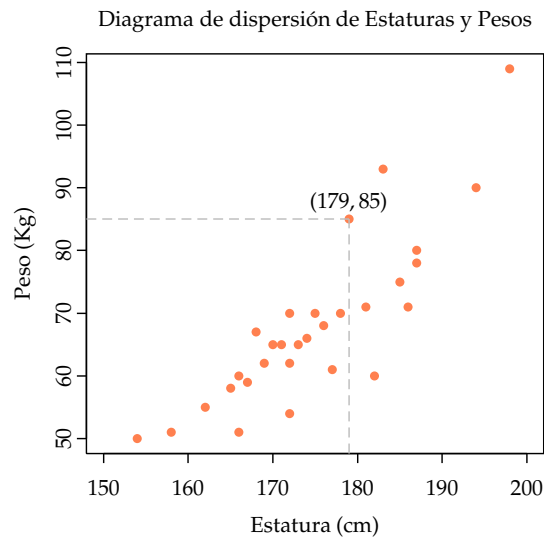


Figura 4.1 – Diagrama de dispersión. El punto (179,85) indicado corresponde a un individuo de la muestra que mide 179 cm y pesa 85 Kg.

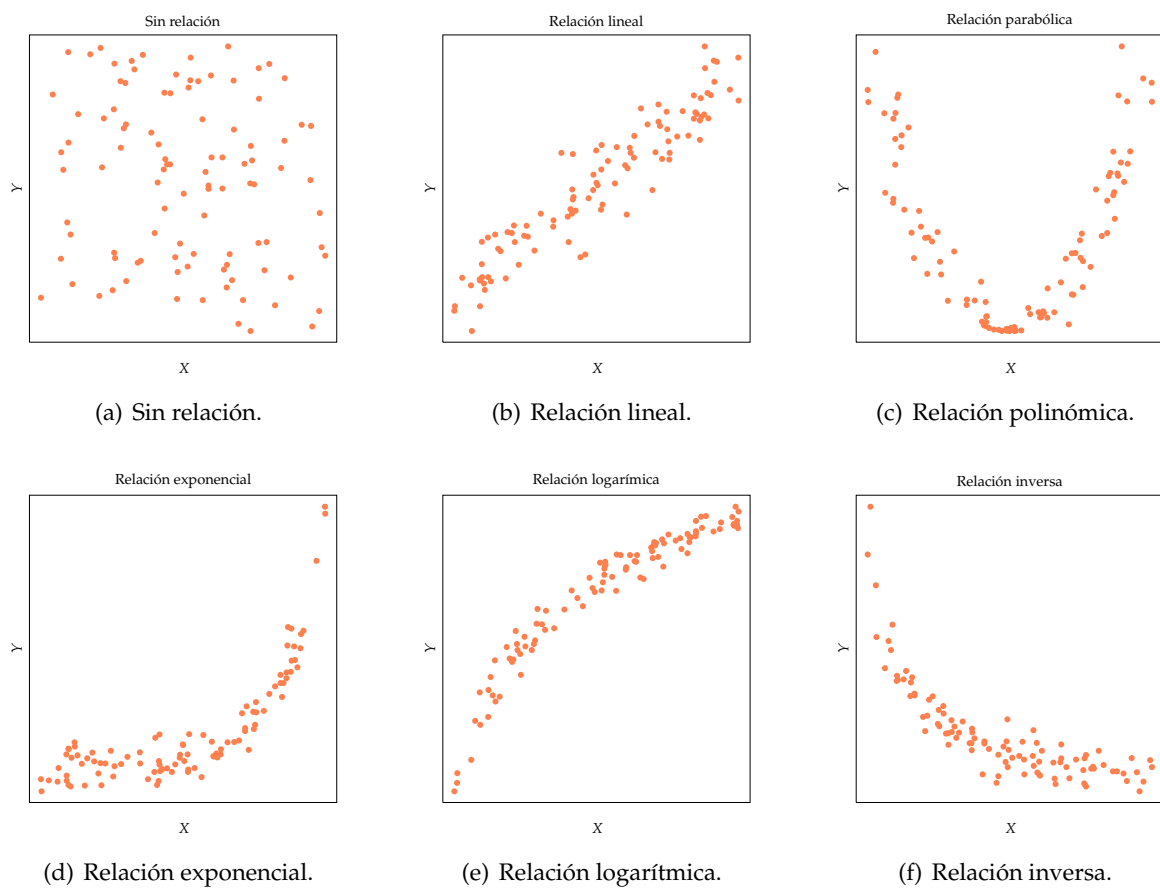


Figura 4.2 – Diagramas de dispersión correspondientes a distintos tipos de relaciones entre variables.

El criterio que suele utilizarse para obtener la función óptima, es que la distancia de cada punto a la curva, medida en el eje Y, sea lo menor posible. A estas distancias se les llama *residuos* o *errores* en Y (figura 4.3). La función que mejor se ajusta a la nube de puntos será, por tanto, aquella que hace mínima la suma de los cuadrados de los residuos.¹

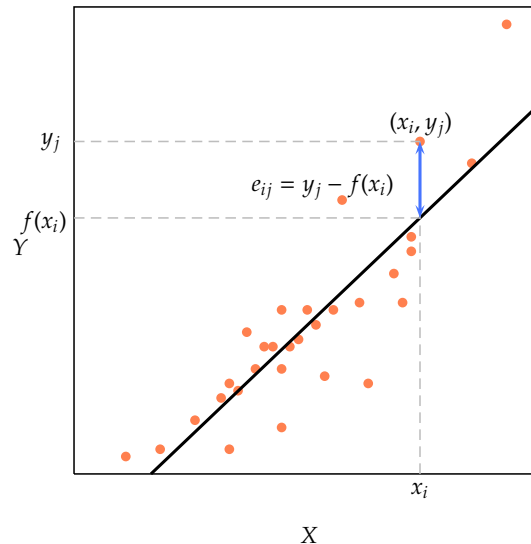


Figura 4.3 – Residuos o errores en Y. El residuo correspondiente a un punto (x_i, y_j) es la diferencia entre el valor y_j observado en la muestra, y el valor teórico del modelo $f(x_i)$, es decir, $e_{ij} = y_j - f(x_i)$.

Rectas de regresión

En el caso de que la nube de puntos tenga forma lineal y optemos por explicar la relación entre X e Y mediante una recta $y = a + bx$, los parámetros a determinar son a (punto de corte con el eje de ordenadas) y b (pendiente de la recta). Los valores de estos parámetros que hacen mínima la suma de residuos al cuadrado, determinan la recta óptima. Esta recta se conoce como *recta de regresión de Y sobre X* y explica la variable Y en función de la variable X. Su ecuación es

$$y = \bar{y} + \frac{s_{xy}}{s_x^2}(x - \bar{x}),$$

donde s_{xy} es un estadístico llamado *covarianza* que mide el grado de relación lineal, y cuya fórmula es

$$s_{xy} = \frac{1}{n} \sum_{i,j} (x_i - \bar{x})(y_j - \bar{y})n_{ij}.$$

■ **Ejemplo 4.2** En la figura 4.4 aparecen las rectas de regresión de Estatura sobre Peso y de Peso sobre Estatura del ejemplo anterior.

La pendiente de la recta de regresión de Y sobre X se conoce como *coeficiente de regresión de Y sobre X*, y mide el incremento que sufrirá la variable Y por cada unidad que se incremente la variable X, según la recta.

Cuanto más pequeños sean los residuos, en valor absoluto, mejor se ajustará el modelo a la nube de puntos, y por tanto, mejor explicará la relación entre X e Y. Cuando todos los residuos son nulos, la recta pasa por todos los puntos de la nube, y la relación es perfecta. En este caso ambas rectas, la de Y sobre X y la de X sobre Y coinciden (figura 4.5(a)).

¹Se elevan al cuadrado para evitar que en la suma se compensen los residuos positivos con los negativos.

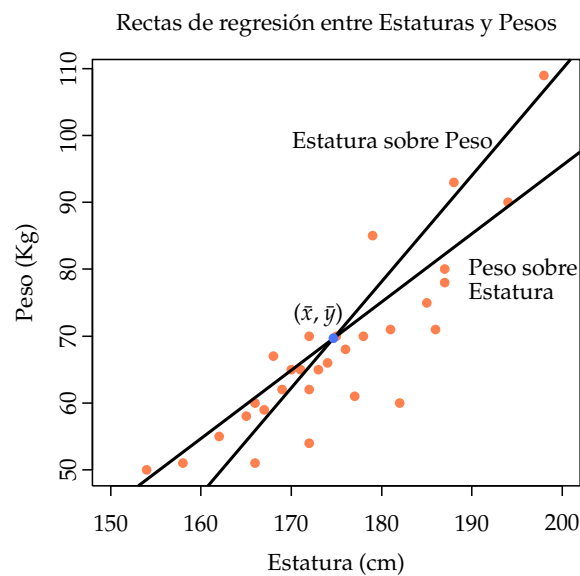
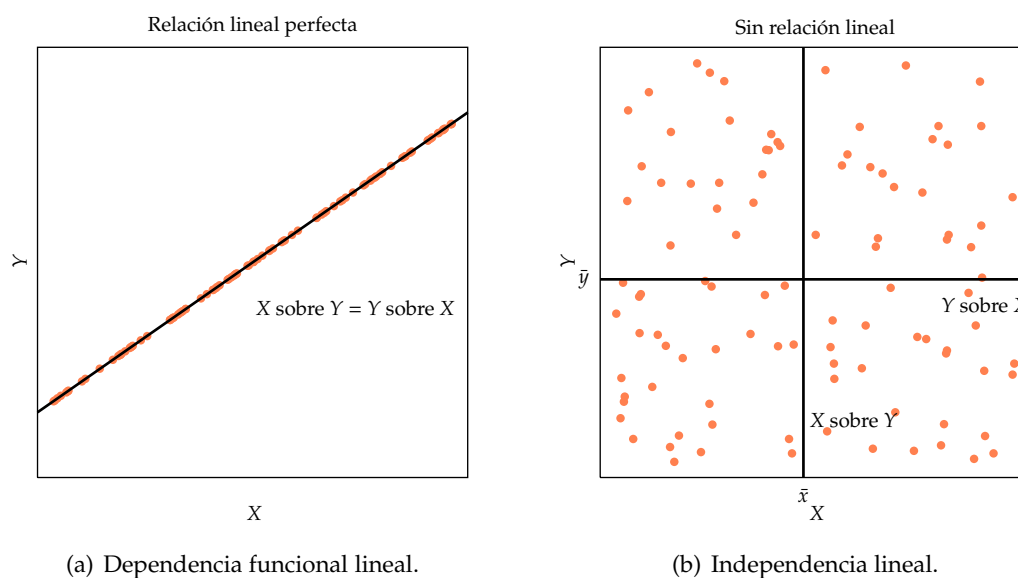


Figura 4.4 – Rectas de regresión de Estatura sobre Peso y de Peso sobre Estatura. Las rectas de regresión siempre se cortan en el punto de medias (\bar{x}, \bar{y})

Por contra, cuando no existe relación lineal entre las variables, la recta de regresión de Y sobre X tiene pendiente nula, y por tanto la ecuación es $y = \bar{y}$, en la que, efectivamente no aparece x , o $x = \bar{x}$ en el caso de la recta de regresión X sobre Y , de manera que ambas rectas se cortan perpendicularmente (figura 4.5(b)).



(a) Dependencia funcional lineal.

(b) Independencia lineal.

Figura 4.5 – Distintos grados de dependencia. En el primer caso, la relación es perfecta y los residuos son nulos. En el segundo caso no existe relación lineal y la pendiente de la recta es nula.

1.2 Correlación

El principal objetivo de la regresión simple es construir un modelo funcional $y = f(x)$ que explique lo mejor posible la relación entre dos variables X (variable independiente) e Y (variable dependiente) medidas en una misma muestra. Generalmente, el modelo construido se utiliza para realizar inferencias predictivas de Y en función de X en el resto de la población. Pero aunque la regresión garantiza que el modelo construido es el mejor posible, dentro del tipo de modelo elegido (lineal, polinómico, exponencial, logarítmico, etc.), puede que aún así, no sea un buen modelo para hacer predicciones, precisamente porque no haya relación de ese tipo entre X e Y . Así pues, con el fin de validar un modelo para realizar predicciones fiables, se necesitan medidas que nos hablen del grado de dependencia entre X e Y , con respecto a un modelo de regresión construido. Estas medidas se conocen como medidas de *correlación*.

Dependiendo del tipo de modelo ajustado, habrá distintos tipos de medidas de correlación. Así, si el modelo de regresión construido es una recta, hablaremos de correlación lineal; si es un polinomio, hablaremos de correlación polinómica; si es una función exponencial, hablaremos de correlación exponencial, etc. En cualquier caso, estas medidas nos hablarán de lo bueno que es el modelo construido, y como consecuencia, de si podemos fiarnos de las predicciones realizadas con dicho modelo.

La mayoría de las medidas de correlación surgen del estudio de los residuos o errores en Y , que son las distancias de los puntos del diagrama de dispersión a la curva de regresión construida, medidas en el eje Y , tal y como se muestra en la figura (4.3). Estas distancias, son en realidad, los errores predictivos del modelo sobre los propios valores de la muestra.

Cuanto más pequeños sean los residuos, mejor se ajustará el modelo a la nube de puntos, y por tanto, mejor explicará la relación entre X e Y . Cuando todos los residuos son nulos, la curva de regresión pasa por todos los puntos de la nube, y entonces se dice que la relación es perfecta, o bien que existe una dependencia funcional entre X e Y (figura 4.5(a)). Por contra, cuando los residuos sean grandes, el modelo no explicará bien la relación entre X e Y , y por tanto, sus predicciones no serán fiables (figura 4.5(b)).

Varianza residual

Una primera medida de correlación, construida a partir de los residuos es la *varianza residual*, que se define como el promedio de los residuos al cuadrado:

$$s_{ry}^2 = \frac{\sum_{i,j} e_{ij}^2 n_{ij}}{n} = \frac{\sum_{i,j} (y_j - f(x_i))^2 n_{ij}}{n}.$$

Cuando los residuos son nulos, entonces $s_{ry}^2 = 0$ y eso indica que hay dependencia funcional. Por otro lado, cuando las variables son independientes, con respecto al modelo de regresión ajustado, entonces los residuos se convierten en las desviaciones de los valores de Y con respecto a su media, y se cumple que $s_{ry}^2 = s_y^2$. Así pues, se cumple que

$$0 \leq s_{ry}^2 \leq s_y^2.$$

Según esto, cuanto menor sea la varianza residual, mayor será la dependencia entre X e Y , de acuerdo al modelo ajustado. No obstante, la varianza tiene como unidades las unidades de Y al cuadrado, y eso dificulta su interpretación.

Coefficiente de determinación

Puesto que el valor máximo que puede tomar la varianza residual es la varianza de Y , se puede definir fácilmente un coeficiente a partir de la comparación de ambas medidas. Surge así el *coeficiente de determinación* que se define como

$$R^2 = 1 - \frac{s_{ry}^2}{s_y^2}.$$

Se cumple que

$$0 \leq R^2 \leq 1,$$

y además no tiene unidades, por lo que es más fácil de interpretar que la varianza residual:

- $R^2 = 0$ indica que existe independencia según el tipo de relación planteada por el modelo de regresión.
- $R^2 = 1$ indica dependencia funcional.

Por tanto, cuanto mayor sea R^2 , mejor será el modelo de regresión.

Si multiplicamos el coeficiente de determinación por 100, se obtiene el porcentaje de variabilidad de Y que explica el modelo de regresión. El porcentaje restante corresponde a la variabilidad que queda por explicar y se corresponde con el error predictivo del modelo. Así, por ejemplo, si tenemos un coeficiente de determinación $R^2 = 0,5$, el modelo de regresión explicaría la mitad de la variabilidad de Y , y en consecuencia, si se utiliza dicho modelo para hacer predicciones, estas tendrían la mitad de error que si no se utilizase, y se tomase como valor de la predicción el valor de la media de Y .

Coeficiente de determinación lineal

En el caso de que el modelo de regresión sea lineal, la fórmula del coeficiente de determinación se simplifica y se convierte en

$$r^2 = \frac{s_{xy}^2}{s_x^2 s_y^2},$$

que se conoce como *coeficiente de determinación lineal*.

Coeficiente de correlación

Otra medida de dependencia bastante habitual es el *coeficiente de correlación*, que se define como la raíz cuadrada del coeficiente de determinación:

$$R = \pm \sqrt{1 - \frac{s_{ry}^2}{s_y^2}},$$

tomando la raíz del mismo signo que la covarianza.

La única ventaja del coeficiente de correlación con respecto al coeficiente de determinación, es que tiene signo, y por tanto, además del grado de dependencia entre X e Y , también nos habla de si la relación es directa (signo +) o inversa (signo -). Su interpretación es:

- $R = 0$ indica independencia con respecto al tipo de relación planteada por el modelo de regresión.
- $R = -1$ indica dependencia funcional inversa.
- $R = 1$ indica dependencia funcional directa.

Por consiguiente, cuanto más próximo esté a -1 o a 1, mejor será el modelo de regresión.

Coeficiente de correlación lineal Al igual que ocurría con el coeficiente de determinación, cuando el modelo de regresión es lineal, la fórmula del coeficiente de correlación se convierte en

$$r = \frac{s_{xy}}{s_x s_y},$$

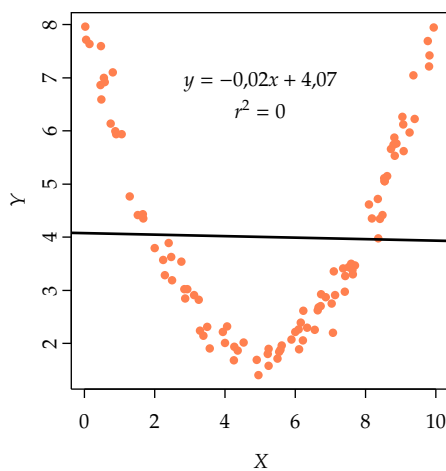
y se llama *coeficiente de correlación lineal*.

Por último, conviene remarcar que un coeficiente de determinación o de correlación nulo, indica que hay independencia según el modelo de regresión construido, pero puede haber dependencia de otro tipo. Esto se ve claramente en el ejemplo de la figura 4.6.

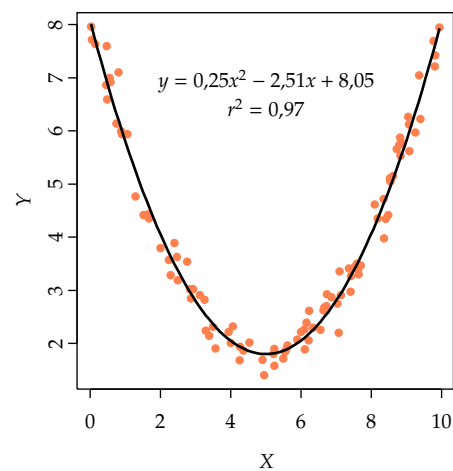
Fiabilidad de las predicciones

Aunque el coeficiente de determinación o de correlación nos hablan de la bondad de un modelo de regresión, no es el único dato que hay que tener en cuenta a la hora de hacer predicciones.

La fiabilidad de las predicciones que hagamos con un modelo de regresión depende de varias cosas:



(a) Dependencia lineal débil.



(b) Dependencia parabólica fuerte.

Figura 4.6 – En la figura de la izquierda se ha ajustado un modelo lineal y se ha obtenido un $R^2 = 0$, lo que indica que el modelo no explica nada de la relación entre X e Y, pero no podemos afirmar que X e Y son independientes. De hecho, en la figura de la derecha se observa que al ajustar un modelo parabólico, $R^2 = 0,97$, lo que indica que casi hay una dependencia funcional parabólica entre X e Y.

- El coeficiente de determinación: Cuando mayor sea, menores serán los errores predictivos y mayor la fiabilidad de las predicciones.
- La variabilidad de la población: Cuanto más variable es una población, más difícil es predecir y por tanto menos fiables serán las predicciones del modelo.
- El tamaño muestral: Cuanto mayor sea, más información tendremos y, en consecuencia, más fiables serán las predicciones.

Además, hay que tener en cuenta que un modelo de regresión es válido para el rango de valores observados en la muestra, pero fuera de ese rango no tenemos información del tipo de relación entre las variables, por lo que no deberíamos hacer predicciones para valores que estén lejos de los observados en la muestra.

2 Ejercicios resueltos

1. Se han medido dos variables X e Y en 10 individuos obteniendo los siguientes resultados:

X	0	1	2	3	4	5	6	7	8	9
Y	2	5	8	11	14	17	20	23	26	29

Se pide:

- Crear un conjunto de datos con las variables X y Y e introducir estos datos.
- Dibujar el diagrama de dispersión correspondiente.



- Seleccionar el menú Teaching►Gráficos►Diagrama de Dispersión.
- En el cuadro de diálogo que aparece, seleccionar la variable Y en el campo Variable Y , la variable X en el campo Variable X , y hacer clic en el botón Enviar.

En vista del diagrama, ¿qué tipo de modelo crees que explicará mejor la relación entre X e Y ?

- Calcular la recta de regresión de Y sobre X .



- Seleccionar el menú Teaching►Regresión►Regresión lineal.
- En el cuadro de diálogo que aparece, seleccionar la variable Y en el campo Variable dependiente y la variable X en el campo Variable independiente, y hacer clic sobre el botón Enviar.

- Dibujar dicha recta sobre el diagrama de dispersión.



- Seleccionar el menú Teaching►Gráficos►Diagrama de Dispersión.
- En el cuadro de diálogo que aparece, seleccionar la variable Y en el campo Variable Y , la variable X en el campo Variable X , y hacer clic en el botón Enviar.
- En la solapa Línea de ajuste, seleccionar Dibujar recta de regresión y hacer clic en el botón Enviar.

- Calcular la recta de regresión de X sobre Y y dibujarla sobre el correspondiente diagrama de dispersión.



Repetir los pasos de los apartados anteriores pero escogiendo como Variable dependiente la variable X , y como Variable independiente la variable Y

- ¿Son grandes los residuos? Comentar los resultados.

2. En una licenciatura se quiere estudiar la relación entre el número medio de horas de estudio diarias y el número de asignaturas suspensas. Para ello se obtuvo la siguiente muestra:

Horas	Suspensos	Horas	Suspensos	Horas	Suspensos
3,5	1	2,2	2	1,3	4
0,6	5	3,3	0	3,1	0
2,8	1	1,7	3	2,3	2
2,5	3	1,1	3	3,2	2
2,6	1	2,0	3	0,9	4
3,9	0	3,5	0	1,7	2
1,5	3	2,1	2	0,2	5
0,7	3	1,8	2	2,9	1
3,6	1	1,1	4	1,0	3
3,7	1	0,7	4	2,3	2

Se pide:

- a) Crear un conjunto de datos con las variables horas.estudio y suspensos e introducir estos datos.
- b) Construir la tabla de frecuencias bidimensional de las variables horas.estudio y suspensos.



- 1) Seleccionar el menú Teaching►Distribución de frecuencias►Tabla de frecuencias bidimensional.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable horas.estudio en el campo Variable a tabular en filas, la variable suspensos en el campo Variable a tabular en columnas, y hacer clic sobre el botón Enviar.

- c) Calcular la recta de regresión de suspensos sobre horas.estudio y dibujarla.



Para calcular la recta de regresión:

- 1) Seleccionar el menú Teaching►Regresión►Regresión lineal.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable suspensos en el campo Variable dependiente y la variable horas.estudio en el campo Variable independiente, seleccionar Guardar el modelo, introducir un nombre para el modelo y hacer clic sobre el botón Enviar.

Para dibujar la recta de regresión:

- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de Dispersión.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable suspensos en el campo Variable Y y la variable horas.estudio en el campo Variable X.
- 3) En la solapa Línea de ajuste, seleccionar Lineal y hacer clic en el botón Enviar.

- d) Indicar el coeficiente de regresión de suspensos sobre horas.estudio. ¿Cómo lo interpretarías?



El coeficiente de regresión es la pendiente de la recta de regresión.

- e) La relación lineal entre estas dos variables, ¿es mejor o peor que la del ejercicio anterior? Comentar los resultados a partir las gráficas de las rectas de regresión y sus residuos.
- f) Calcular los coeficientes de correlación y de determinación lineal. ¿Es un buen modelo la recta de regresión? ¿Qué porcentaje de la variabilidad del número de suspensos está explicada por el modelo?



El coeficiente de determinación aparece en la ventana de resultados como R^2 ajustado, y el coeficiente de correlación es su raíz cuadrada.

- g) Utilizar la recta de regresión para predecir el número de suspensos correspondiente a 3 horas de estudio diarias. ¿Es fiable esta predicción?



- 1) Seleccionar el menú Teaching►Regresión►Predicciones.
- 2) En el cuadro de diálogo que aparece seleccionar como modelo de regresión la recta calculada en el segundo apartado, introducir los valores para los que se desea la predicción en el campo Predicciones para y hacer clic sobre el botón Enviar.

- h) Según el modelo lineal, ¿cuántas horas diarias tendrá que estudiar como mínimo un alumno si quiere aprobarlo todo?



Seguir los mismos pasos de los apartados anteriores, pero escogiendo como variable dependiente `horas.estudio`, y como independiente `suspensos`, y haciendo la predicción para 0 `suspensos`.

3. Después de tomar un litro de vino se ha medido la concentración de alcohol en la sangre en distintos instantes, obteniendo:

Tiempo después (minutos)	30	60	90	120	150	180	210
Concentración (gramos/litro)	1,6	1,7	1,5	1,1	0,7	0,2	2,1

Se pide:

- Crear las variables tiempo y alcohol e introducir estos datos.
- Calcular el coeficiente de correlación lineal entre el alcohol y el tiempo e interpretarlo. ¿Es bueno el modelo lineal?



- Seleccionar el menú Teaching►Regresión►Regresión lineal.
- En el cuadro de diálogo que aparece, seleccionar la variable alcohol en el campo Variable dependiente y la variable tiempo en el campo Variable independiente, y hacer clic sobre el botón Enviar.

- Dibujar la recta de regresión del alcohol sobre el tiempo. ¿Existe algún individuo con un residuo demasiado grande? Si es así, eliminar dicho individuo de la muestra y volver a calcular el coeficiente de correlación. ¿Ha mejorado el modelo?



- Seleccionar el menú Teaching►Gráficos►Diagrama de Dispersión.
- En el cuadro de diálogo que aparece, seleccionar la variable alcohol en el campo Variable Y y la variable tiempo en el campo Variable X.
- En la solapa Línea de ajuste, seleccionar Lineal y hacer clic en el botón Enviar.

Se observa que hay un residuo atípico para el punto que corresponde a los 210 minutos. Para eliminarlo: En la ventana de edición del conjunto de datos hacer clic con el botón derecho del ratón sobre la fila correspondiente al dato con el residuo atípico y seleccionar Borrar esta fila.

- Si la concentración máxima de alcohol en la sangre que permite la ley para poder conducir es 0,3 g/l, ¿cuánto tiempo habrá que esperar después de tomarse un litro de vino para poder conducir sin infringir la ley? ¿Es fiable esta predicción?



Para construir la recta de regresión:

- Seleccionar el menú Teaching►Regresión►Regresión lineal.
- En el cuadro de diálogo que aparece, seleccionar la variable tiempo en el campo Variable dependiente y la variable alcohol en el campo Variable independiente.
- Seleccionar Guardar el modelo, introducir un nombre para el modelo y hacer clic sobre el botón Enviar.

Para hacer la predicción:

- Seleccionar el menú Teaching►Regresión►Predicciones.
- En el cuadro de diálogo que aparece seleccionar como modelo de regresión la recta calculada e introducir los valores para los que se desea la predicción en el campo Predicciones para y hacer clic sobre el botón Enviar.

4. El conjunto de datos `edad.estatura` del paquete `rk.Teaching` contiene la edad y la estatura de 30 personas. Se pide:

4. Regresión Lineal Simple y Correlación

- a) Cargar datos del conjunto de datos edad.estatura desde el paquete rk.Teaching.
- b) Calcular la recta de regresión de la estatura sobre la edad. ¿Es un buen modelo la recta de regresión?



- 1) Seleccionar el menú Teaching►Regresión►Regresión lineal.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable estatura en el campo Variable dependiente y la variable edad en el campo Variable independiente, y hacer clic en el botón Enviar.

- c) Dibujar el diagrama de dispersión de la estatura sobre la edad. ¿Alrededor de qué edad se observa un cambio en la tendencia?



- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de Dispersión.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable estatura en el campo Variable Y, la variable edad en el campo Variable X, y hacer clic en el botón Enviar.

- d) Recodificar la variable edad en dos grupos para mayores y menores de 20 años.



- 1) Seleccionar el menú Teaching►Datos►Recodificar variable.
- 2) En el cuadro de diálogo que aparece seleccionar en el campo Variable a recodificar la variable edad.
- 3) En el campo Reglas de recodificación introducir
10:20 = "menores"
20:hi = "mayores"
- 4) En el cuadro Guardar nueva variable hacer clic sobre el botón Cambiar.
- 5) En el cuadro de diálogo que aparece seleccionar como objeto padre la el conjunto de datos edad_estatura y hacer clic sobre el botón Aceptar.
- 6) Introducir el nombre de la nueva variable grupo.edad y hacer clic sobre el botón Enviar.

- e) Calcular la recta de regresión de la estatura sobre la edad para cada grupo de edad. ¿En qué grupo explica mejor la recta de regresión la relación entre la estatura y la edad? Justificar la respuesta.



- 1) Seleccionar el menú Teaching►Regresión►Regresión lineal.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable estatura en el campo Variable dependiente y la variable edad como Variable independiente.
- 3) Seleccionar la opción Ajuste por grupos, introducir la variable grupo.edad en el campo Variable de agrupación, y hacer clic en el Enviar.

- f) Dibujar las rectas de regresión anteriores.



- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de Dispersión.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable estatura en el campo Variable Y y la variable edad en el campo Variable X.
- 3) Seleccionar la opción Dibujar por grupos e introducir la variable grupo.edad en el campo Variable de agrupación.
- 4) En la solapa Línea de ajuste, seleccionar Lineal y hacer clic en el botón Enviar.

- g) ¿Qué estatura se espera que tenga una persona de 14 años? ¿Y una de 38?



Para predecir la estatura de la persona de 14 años:

- 1) Seleccionar el menú Teaching►Regresión►Predicciones.
- 2) En el cuadro de diálogo que aparece seleccionar como modelo de regresión la recta calculada para los menores e introducir 14 en el campo Predicciones para y hacer clic sobre el botón Enviar.

para predecir la estatura de la persona de 38 años, repetir lo mismo pero seleccionando la recta de regresión para los mayores e introducir 38 en el campo Predicciones para.

5. La siguiente tabla recoge la información de las calificaciones obtenidas por un grupo de alumnos en dos asignaturas X e Y.

Alumno	1	2	3	4	5	6	7	8	9	10	11	12
X	NT	AP	SS	SS	AP	AP	SS	NT	SB	SS	AP	AP
Y	SB	SS	AP	SS	AP	NT	SS	NT	NT	AP	AP	NT

Se pide:

- a) Crear un conjunto de datos con las variables X e Y e introducir los datos.
- b) ¿Existe relación entre las calificaciones de X e Y? Justificar la respuesta.



- 1) Seleccionar el menú Teaching►Regresión►Correlación.
- 2) En el cuadro de diálogo que aparece seleccionar la variables X e Y en el campo Variables.
- 3) En la solapa Opciones de correlación seleccionar el método de Ro de Spearman y hacer clic sobre el botón Enviar.

3 Ejercicios propuestos

1. Se determina la pérdida de actividad que experimenta un medicamento desde el momento de su fabricación a lo largo del tiempo, obteniéndose el siguiente resultado:

Tiempo (en años)	1	2	3	4	5
Actividad restante (%)	96	84	70	58	52

Se desea calcular:

- a) La relación fundamental (recta de regresión) entre actividad restante y tiempo transcurrido.
 - b) ¿En qué porcentaje disminuye la actividad cada año que pasa?
 - c) ¿Cuándo tiempo debe pasar para que el fármaco tenga una actividad del 80 %? ¿Cuándo será nula la actividad? ¿Son igualmente fiables estas predicciones?
2. Al realizar un estudio sobre la dosificación de un cierto medicamento, se trataron 6 pacientes con dosis diarias de 2 mg, 7 pacientes con 3 mg y otros 7 pacientes con 4 mg. De los pacientes tratados con 2 mg, 2 curaron al cabo de 5 días, y 4 al cabo de 6 días. De los pacientes tratados con 3 mg diarios, 2 curaron al cabo de 3 días, 4 al cabo de 5 días y 1 al cabo de 6 días. Y de los pacientes tratados con 4 mg diarios, 5 curaron al cabo de 3 días y 2 al cabo de 4 días. Se pide:
- a) Calcular la recta de regresión del tiempo de curación con respecto a la dosis suministrada.
 - b) Calcular el coeficiente de regresión del tiempo de curación con respecto a la dosis e interpretarlo.
 - c) Calcular el coeficiente de correlación lineal e interpretarlo.

4. Regresión Lineal Simple y Correlación

- d) Determinar el tiempo esperado de curación para una dosis de 5 mg diarios. ¿Es fiable esta predicción?
 - e) ¿Qué dosis debe aplicarse si queremos que el paciente tarde 4 días en curarse? ¿Es fiable la predicción?
3. El fichero estaturas.pesos.alumnos del paquete rk.Teaching, contiene la estatura, el peso y el sexo de una muestra de alumnos universitarios. Se pide:
- a) Cargar el conjunto de datos estaturas.pesos.alumnos desde el paquete rk.Teaching.
 - b) Calcular la recta de regresión del peso sobre la estatura y dibujarla.
 - c) Calcular las rectas de regresión del peso sobre la estatura para cada sexo y dibujarlas.
 - d) Calcular los coeficientes de determinación de ambas rectas. ¿Qué recta es mejor modelo? Justificar la respuesta.
 - e) ¿Qué peso tendrá un hombre que mida 170 cm? ¿Y una mujer de la misma estatura?
4. El conjunto de datos neonatos del paquete rk.Teaching, contiene información sobre una muestra de 320 recién nacidos en un hospital durante un año que cumplieron el tiempo normal de gestación. Se pide:
- a) Construir la tabla de frecuencias bidimensional del Agpar al minuto de nacer frente a si la madre ha fumado o no durante el embarazo. ¿Qué conclusiones se pueden sacar?
 - b) Construir la tabla de frecuencias bidimensional del peso de los recién nacidos frente a la edad de la madre. ¿Qué conclusiones se pueden sacar?
 - c) Construir la recta de regresión del peso de los recién nacidos sobre el número de cigarros fumados al día por las madres. ¿Existe una relación lineal fuerte entre el peso y el número de cigarros?
 - d) Dibujar la recta de regresión calculada en el apartado anterior. ¿Por qué la recta no se ajusta bien a la nube de puntos?
 - e) Calcular y dibujar la recta de regresión del peso de los recién nacidos sobre el número de cigarros fumados al día por las madres en el grupo de las madres que si fumaron durante el embarazo. ¿Es este modelo mejor o peor que la recta de los apartados anteriores?
Según este modelo, ¿cuánto disminuirá el peso del recién nacido por cada cigarro más diario que fume la madre?
 - f) Según el modelo anterior, ¿qué peso tendrá un recién nacido de una madre que ha fumado 5 cigarros diarios durante el embarazo? ¿Y si la madre ha fumado 30 cigarros diarios durante el embarazo? ¿Son fiables estas predicciones?
 - g) ¿Existe la misma relación lineal entre el peso de los recién nacidos y el número de cigarros fumados al día por las madres que fumaron durante el embarazo en el grupo de las madres menores de 20 y en el grupo de las madres mayores de 20? ¿Qué se puede concluir?

Regresión no lineal

1 Fundamentos teóricos

La regresión simple tiene por objeto la construcción de un modelo funcional $y = f(x)$ que explique lo mejor posible la relación entre dos variables Y (variable dependiente) y X (variable independiente) medidas en una misma muestra.

Ya vimos que, dependiendo de la forma de esta función, existen muchos tipos de regresión simple. Entre los más habituales están:

Modelo	Ecuación genérica
Lineal	$y = a + bx$
Parabólico	$y = a + bx + cx^2$
Polinómico de grado n	$y = a_0 + a_1x + \dots + a_nx^n$
Potencial	$y = ax^b$
Exponencial	$y = e^{a+bx}$
Logarítmico	$y = a + b \log x$
Inverso	$y = a + b/x$
Curva S	$y = e^{a+b/x}$

La elección de un tipo de modelo u otro suele hacerse según la forma de la nube de puntos del diagrama de dispersión. A veces estará claro qué tipo de modelo se debe construir, tal y como ocurre en los diagramas de dispersión de la figura 5.1. Pero otras veces no estará tan claro, y en estas ocasiones, lo normal es ajustar los dos o tres modelos que nos parezcan más convincentes, para luego quedarnos con el que mejor explique la relación entre Y y X , mirando el coeficiente de determinación¹ de cada modelo.

Ya vimos en la práctica sobre regresión lineal simple, cómo construir rectas de regresión. En el caso de que optemos por ajustar un modelo no lineal, la construcción del mismo puede realizarse siguiendo los mismos pasos que en el caso lineal. Básicamente se trata de determinar los parámetros del modelo que minimizan la suma de los cuadrados de los residuos en Y . En los modelos multiplicativo y exponencial, el sistema aplica transformaciones logarítmicas a las variables y después ajusta un modelo lineal a los datos transformados. En el modelo recíproco, el sistema sustituye la variable dependiente por su recíproco antes de estimar la ecuación de regresión.

¹Ver la práctica de regresión lineal y correlación.

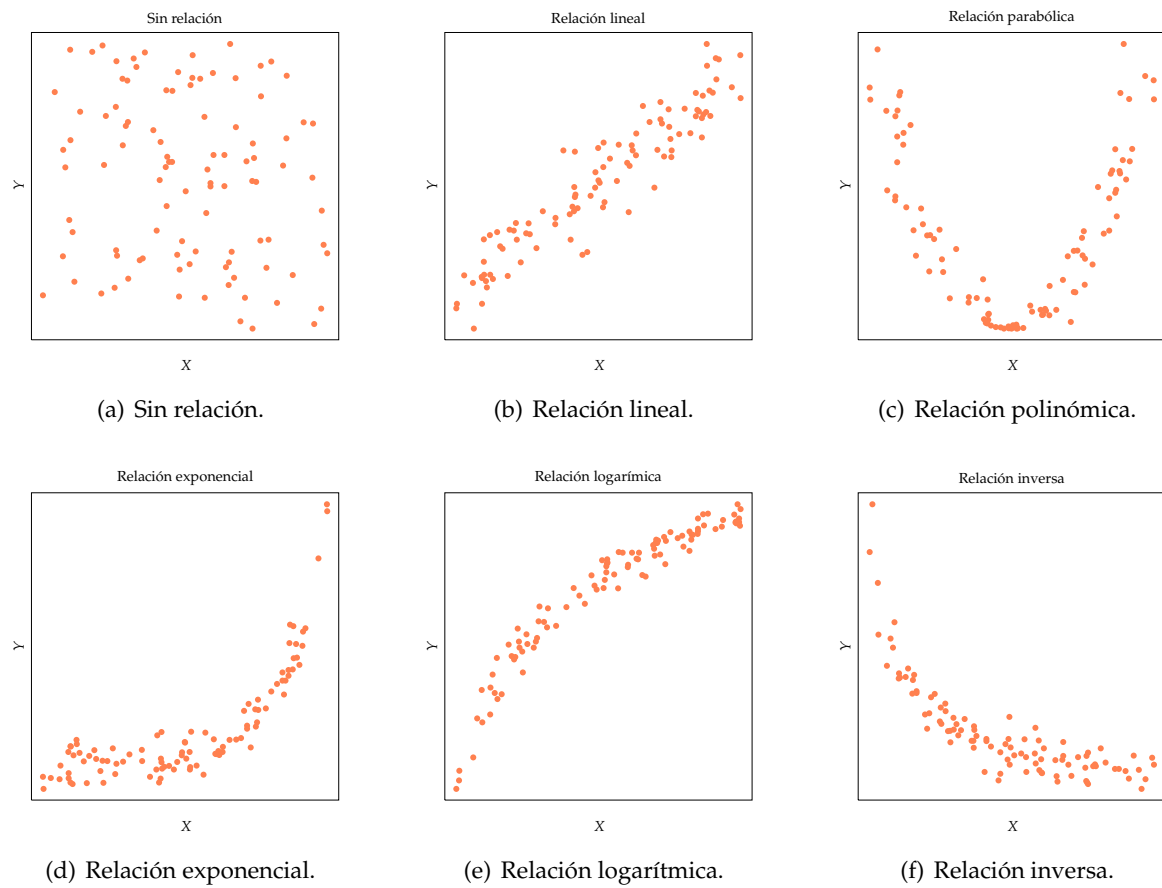


Figura 5.1 – Diagramas de dispersión correspondientes a distintos tipos de relaciones entre variables.

2 Ejercicios resueltos

El procedimiento más sencillo para construir un modelo no lineal, siempre que sea posible, es transformar las variables para convertirlo en un modelo lineal. En el caso de los modelos de regresión simple más comunes las transformaciones que convierten cada modelo en un modelo lineal aparecen en la tabla siguiente:

Modelo	Modelo no lineal	Modelo lineal	Transformación
Potencial	$y = ax^b$	$\log(y) = \log(a) + b \log(x)$	Se toma el logaritmo de ambas variables
Exponencial	$y = e^{a+bx}$	$\log(y) = a + bx$	Se toma el logaritmo de la variable dependiente
Logarítmico	$y = a + b \log x$	$y = a + b \log x$	Se toma el logaritmo de la variable independiente
Inverso	$y = a + b/x$	$y = a + b \frac{1}{x}$	Se toma el inverso de la variable independiente
Curva S	$y = e^{a+b/x}$	$\log(y) = a + b \frac{1}{x}$	Se toma el logaritmo de la variable dependiente y el inverso de la independiente

1. En un experimento se ha medido el número de bacterias por unidad de volumen en un cultivo, cada hora transcurrida, obteniendo los siguientes resultados:

Horas	0	1	2	3	4	5	6	7	8
Nº Bacterias	25	28	47	65	86	121	190	290	362

Se pide:

- a) Crear un conjunto de datos con las variables horas y bacterias e introducir estos datos.
- b) Dibujar el diagrama de dispersión correspondiente. En vista del diagrama, ¿qué tipo de modelo crees que explicará mejor la relación entre el número de bacterias y el tiempo transcurrido?



- 1) Seleccionar el menú Teaching ▶ Gráficos ▶ Diagrama de dispersión.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable bacterias en el campo Variable Y y la variable horas en el campo Variable X, y hacer clic en el botón Enviar.

- c) Calcular los modelos exponencial y cuadrático de las bacterias sobre las horas. ¿Qué tipo de modelo es el mejor?



Para el modelo exponencial:

- 1) Seleccionar el menú Teaching ▶ Regresión ▶ Regresión no lineal.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable bacterias en el campo Variable dependiente y la variable horas en el campo Variable independiente.
- 3) En la solapa de Modelo de regresión seleccionar el modelo Exponencial.
- 4) Seleccionar Guardar modelo e introducir un nombre para el modelo y hacer clic sobre el botón Enviar.

Para el modelo cuadrático repetir los pasos pero seleccionando como modelo el Cuadrático. El modelo mejor será aquel que tenga un coeficiente de determinación mayor.

- d) Dibujar la curva del mejor de los modelos anteriores.



- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de dispersión.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable bacterias en el campo Variable Y y la variable horas en el campo Variable X.
- 3) En la solapa Línea de ajuste seleccionar la opción Exponencial y hacer clic sobre el botón Enviar.

e) Según el modelo anterior, ¿cuántas bacterias habrá al cabo de 3 horas y media del inicio del cultivo? ¿Y al cabo de 10 horas? ¿Son fiables estas predicciones?



- 1) Seleccionar el menú Teaching►Regresión►Predicciones.
- 2) En el cuadro de diálogo que aparece seleccionar el modelo de regresión exponencial construido antes.
- 3) Introducir los valores 3,5, 10 en el campo Predicciones para y hacer clic sobre el botón Enviar.
- 4) Como se trata de un modelo exponencial, las predicciones obtenidas corresponden al logaritmo de bacterias. Para obtener la predicción de bacterias basta con aplicar la función exponencial a los valores obtenidos.

f) Dar una predicción lo más fiable posible del tiempo que tendría que transcurrir para que en el cultivo hubiese 100 bacterias.



Para construir el modelo logarítmico:

- 1) Seleccionar el menú Teaching►Regresión►Regresión no lineal.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable horas en el campo Variable dependiente y la variable bacterias en el campo Variable independiente.
- 3) Seleccionar como modelo el Logarítmico.
- 4) Seleccionar Guardar modelo e introducir un nombre para el modelo y hacer clic sobre el botón Enviar.

Para hacer la predicción:

- 1) Seleccionar el menú Teaching►Regresión►Predicciones.
- 2) En el cuadro de diálogo que aparece seleccionar el modelo de regresión logarítmico construido antes.
- 3) Introducir el valor 100 en el campo Predicciones para y hacer clic sobre el botón Enviar.

2. El conjunto de datos dieta del paquete rk.Teaching contiene los datos de un estudio llevado a cabo por un centro dietético para probar una nueva dieta de adelgazamiento. Para cada individuo se ha medido el número de días que lleva con la dieta, el número de kilos perdidos desde entonces y si realizó o no un programa de ejercicios. Se pide:

- a) Cargar el conjunto de datos dieta desde el paquete rk.Teaching.
- b) Dibujar el diagrama de dispersión. Según la nube de puntos, ¿qué tipo de modelo explicaría mejor la relación entre los kilos perdidos y los días de dieta?



- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de dispersión.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable kilos en el campo Variable Y, la variable días en el campo Variable X, y hacer clic en el botón Enviar.

c) Construir el modelo de regresión que mejor explique la relación entre los kilos perdidos y los días de dieta.



- 1) Seleccionar el menú Teaching►Regresión►Comparación de modelos.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable kilos en el campo Variable dependiente y la variable días en el campo Variable independiente.
- 3) En la solapa Modelos de regresión seleccionar todos los modelos y hacer clic sobre el botón Enviar.
- 4) El mejor modelo aparece en primer lugar y es el que tenga el coeficiente de determinación mayor.

d) Dibujar el modelo del apartado anterior.



- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de Dispersión.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable kilos en el campo Variable Y y la variable días en el campo Variable X.
- 3) En la solapa Línea de ajuste seleccionar la opción correspondiente al mejor modelo y hacer clic sobre el botón Enviar.

e) Construir el modelo de regresión que mejor explique la relación entre los kilos perdidos y los días de dieta para los que no hacen ejercicio.



Para ver qué modelo es mejor:

- 1) Seleccionar el menú Teaching►Regresión►Comparación de modelos.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable kilos en el campo Variable dependiente y la variable días en el campo Variable independiente.
- 3) Seleccionar la opción Filtro e introducir la condición ejercicio=="no" en el campo Condición de selección.
- 4) En la solapa Modelos de regresión seleccionar todos los modelos y hacer clic sobre el botón Enviar.
- 5) El mejor modelo aparece en primer lugar y es el que tenga el coeficiente de determinación mayor.

Para construir el modelo:

- 1) Seleccionar el menú Teaching►Regresión►Regresión no lineal.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable kilos en el campo Variable dependiente y la variable días en el campo Variable independiente.
- 3) Seleccionar la opción Filtro e introducir la condición ejercicio=="no" en el campo Condición de selección.
- 4) Seleccionar Guardar modelo e introducir un nombre para el modelo y hacer clic sobre el botón Enviar.

f) Construir el modelo de regresión que mejor explique la relación entre los kilos perdidos y los días de dieta para los que si hacen ejercicio.



Para ver qué modelo es mejor:

- 1) Seleccionar el menú Teaching►Regresión►Comparación de modelos.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable kilos en el campo Variable dependiente y la variable días en el campo Variable independiente.
- 3) Seleccionar la opción Filtro e introducir la condición ejercicio=="si" en el campo Condición de selección.
- 4) En la solapa Modelos de regresión seleccionar todos los modelos y hacer clic sobre el botón Enviar.

- 5) El mejor modelo aparece en primer lugar y es el que tenga el coeficiente de determinación mayor.

Para construir el modelo:

- 1) Seleccionar el menú Teaching ▶ Regresión ▶ Regresión no lineal.
- 2) En el cuadro de diálogo que aparece, seleccionar la variable kilos en el campo Variable dependiente y la variable dias en el campo Variable independiente.
- 3) Seleccionar la opción Filtro e introducir la condición ejercicio=="si" en el campo Condición de selección.
- 4) Seleccionar Guardar modelo e introducir un nombre para el modelo y hacer clic sobre el botón Enviar.

- g) Utilizar el modelo construido para predecir el número de kilos perdidos tras 40 y 500 días de dieta, tanto para los que hacen ejercicio como para los que no. ¿Son fiables estas predicciones?



- 1) Seleccionar el menú Teaching ▶ Regresión ▶ Predicciones.
- 2) En el cuadro de diálogo que aparece seleccionar el modelo de regresión construido antes para los que no hacen ejercicio.
- 3) Introducir los valores 40, 500 en el campo Predicciones para y hacer clic sobre el botón Enviar.

Repetir los pasos anteriores seleccionando el modelo de regresión construido antes para los que si hacen ejercicio.

3 Ejercicios propuestos

1. La concentración de un fármaco en sangre, C en mg/dl, es función del tiempo, t en horas, y viene dada por la siguiente tabla:

t	2	3	4	5	6	7	8
C	25	36	48	64	86	114	168

Se pide:

- a) Según el modelo exponencial, ¿qué concentración de fármaco habría a las 4,8 horas? ¿Es fiable la predicción? Justificar adecuadamente la respuesta.
- b) Según el modelo logarítmico, ¿qué tiempo debe pasar para que la concentración sea de 100 mg/dl?
2. El fichero naciones.txt contiene información sobre el desarrollo de distintos países (tasa de fertilidad, tasa de uso de anticonceptivos, tasa de mortalidad infantil, producto interior bruto per cápita y continente). Se pide:
 - a) Importar el fichero naciones.txt en un conjunto de datos.
 - b) Construir el mejor modelo de regresión de la tasa de fertilidad sobre el producto interior bruto. ¿Cómo explicarías esta relación?
 - c) Dibujar el modelo del apartado anterior.
 - d) ¿Qué tasa de fertilidad le corresponde a una mujer que viva en un país con un producto interior bruto per cápita de 10000 \$? ¿Y si la mujer vive en Europa?

Probabilidad

1 Fundamentos teóricos

1.1 Introducción

La estadística descriptiva permite describir el comportamiento y las relaciones entre las variables en la muestra, pero no permite sacar conclusiones sobre el resto de la población.

Ha llegado el momento de dar el salto de la muestra a la población y pasar de la estadística descriptiva a la inferencia estadística, y el puente que lo permite es la *teoría de la probabilidad*.

Hay que tener en cuenta que el conocimiento que se puede obtener de la población a partir de la muestra es limitado, pero resulta evidente que la aproximación a la realidad de la población será mejor cuanto más representativa sea la muestra de ésta. Y recordemos que para que la muestra sea representativa de la población deben utilizarse técnicas de muestreo aleatorio, es decir, en la que los individuos se seleccionen al azar.

La teoría de la probabilidad precisamente se encarga de controlar ese azar para saber hasta qué punto son fiables y extrapolables al resto de la población las conclusiones obtenidas a partir de una muestra.

1.2 Experimentos y sucesos aleatorios

El estudio de una característica en una población se realiza a través de experimentos aleatorios.

Definición 6.1 — Experimento aleatorio. Un *experimento aleatorio* es aquel en el que se conoce cuál es el conjunto de resultados posibles antes de su realización pero se desconoce cuál será el resultado concreto del mismo.

Un ejemplo sencillo de experimentos aleatorios son los juegos de azar. Por ejemplo, el lanzamiento de un dado es un experimento aleatorio ya que:

- Se conoce el conjunto posibles de resultados $\{1, 2, 3, 4, 5, 6\}$.
- Antes de lanzar el dado, es imposible predecir con absoluta certeza el valor que saldrá.

Otro ejemplo de experimento aleatorio sería la selección de un individuo de una población al azar y la determinación de su grupo sanguíneo.

En general, la obtención de cualquier muestra mediante procedimientos aleatorios será un experimento aleatorio.

Definición 6.2 — Espacio muestral. Al conjunto E de todos los posibles resultados de un experimento aleatorio se le llama *espacio muestral*.

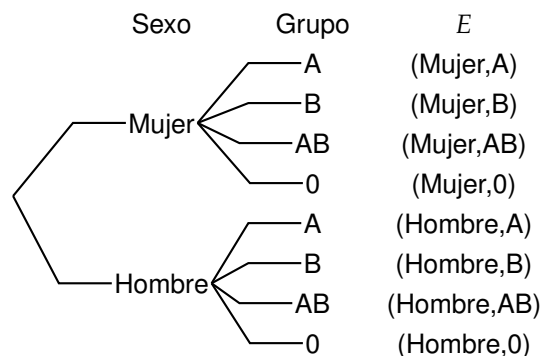
Algunos ejemplos de espacios muestrales son:

- Lanzamiento de una moneda: $E = \{c, x\}$.
- Lanzamiento de un dado: $E = \{1, 2, 3, 4, 5, 6\}$.

- Grupo sanguíneo de un individuo seleccionado al azar: $E = \{A, B, AB, 0\}$.
- Estatura de un individuo seleccionado al azar: \mathbb{R}^+ .

En los experimentos donde se miden más de una variable, la construcción del espacio muestral puede complicarse. En tales casos, es recomendable utilizar un *diagrama de árbol* de manera que cada nivel del árbol es una variable observada y cada rama un posible valor.

Por ejemplo, si el experimento consiste en observar el sexo y el grupo sanguíneo de una persona, el espacio muestral podría construirse mediante el siguiente árbol:



En RKWard los espacios muestrales se representan mediante conjuntos de datos con las variables que se midan en el experimento, indicando en cada fila un resultado posible. Por ejemplo, el conjunto de datos correspondiente al espacio muestral del experimento anterior se muestra en la figura 6.1.

	1	2	#Nueva variable#
Nombre	sexo	grupo.sanguineo	
Etiqueta			
Tipo	Factor	Factor	
Formato			
Niveles	mujer#,#hom...	0#,#A#,#B#,#AB	
1	mujer	A	
2	mujer	B	
3	mujer	AB	
4	mujer	0	
5	hombre	A	
6	hombre	B	
7	hombre	AB	
8	hombre	0	

Figura 6.1 – Conjunto de datos correspondiente al espacio muestral del experimento consistente en sacar un individuo al azar de una población y medir su sexo y su grupo sanguíneo.

Definición 6.3 — Suceso aleatorio. Un *suceso aleatorio* es cualquier subconjunto del espacio muestral E de un experimento aleatorio.

Existen distintos tipos de sucesos:

Suceso imposible: Es el subconjunto vacío \emptyset . El suceso nunca ocurre.

Sucesos elementales: Son los subconjuntos formados por un solo elemento.

Sucesos compuestos: Son los subconjuntos formados por dos o más elementos.

Suceso seguro: Es el propio espacio muestral. El suceso seguro siempre ocurre.

Definición 6.4 — Espacio de sucesos. Dado un espacio muestral E de un experimento aleatorio, el conjunto formado por todos los posibles sucesos de E se llama *espacio de sucesos de E* y se denota $\mathcal{P}(E)$.

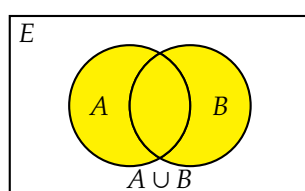
■ **Ejemplo 6.1** Dado el espacio muestral $E = \{a, b, c\}$, se tiene

$$\mathcal{P}(E) = \{\emptyset, \{a\}, \{b\}, \{c\}, \{a, b\}, \{a, c\}, \{b, c\}, \{a, b, c\}\}$$

Puesto que los sucesos son conjuntos, tiene sentido definir operaciones entre sucesos a partir de la teoría de conjuntos.

Definición 6.5 — Suceso unión. Dados dos sucesos $A, B \in \mathcal{P}(E)$, se llama *suceso unión* de A y B , y se denota $A \cup B$, al suceso formado por los elementos de A junto a los elementos de B , es decir,

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

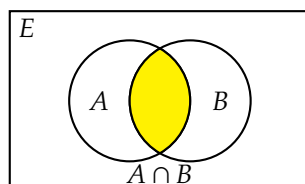


El suceso unión $A \cup B$ ocurre siempre que ocurre A o B .

■ **Ejemplo 6.2** Sea $E = \{1, 2, 3, 4, 5, 6\}$, el conjunto de los números de un dado, y $A = \{2, 4, 6\}$ y $B = \{1, 2, 3, 4\}$. Entonces $A \cup B = \{1, 2, 3, 4, 6\}$.

Definición 6.6 — Suceso intersección. Dados dos sucesos $A, B \in \mathcal{P}(E)$, se llama *suceso intersección* de A y B , y se denota $A \cap B$, al suceso formado por los elementos comunes de A y B , es decir,

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



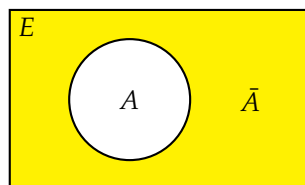
El suceso intersección $A \cap B$ ocurre siempre que ocurren A y B .

■ **Ejemplo 6.3** Sea $E = \{1, 2, 3, 4, 5, 6\}$, el conjunto de los números de un dado, y $A = \{2, 4, 6\}$ y $B = \{1, 2, 3, 4\}$. Entonces $A \cap B = \{2, 4\}$.

Diremos que dos sucesos son **incompatibles** si su intersección es vacía. Por ejemplo $A = \{2, 4, 6\}$ y $C = \{1, 3\}$ son incompatibles.

Definición 6.7 — Suceso contrario. Dado un conjunto $A \in \mathcal{P}(E)$, se llama *suceso contrario o complementario* de A , y se denota \bar{A} , al suceso formado por los elementos de E que no pertenecen a A , es decir,

$$\bar{A} = \{x \mid x \notin A\}.$$

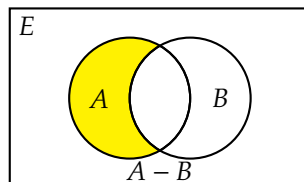


El suceso contrario \bar{A} ocurre siempre que no ocurre A .

■ **Ejemplo 6.4** Sea $E = \{1, 2, 3, 4, 5, 6\}$, el conjunto de los números de un dado, y $A = \{2, 4, 6\}$. Entonces $\bar{A} = \{1, 3, 5\}$.

Definición 6.8 — Suceso diferencia. Dados dos sucesos $A, B \in \mathcal{P}(E)$, se llama *suceso diferencia* de A y B , y se denota $A - B$, al suceso formado por los elementos de A que no pertenecen a B , es decir,

$$A - B = \{x \mid x \in A \text{ y } x \notin B\}.$$



El suceso diferencia $A - B$ ocurre siempre que ocurre A pero no ocurre B , y también puede expresarse como $A \cap \bar{B}$.

■ **Ejemplo 6.5** Sea $E = \{1, 2, 3, 4, 5, 6\}$, el conjunto de los números de un dado, y $A = \{2, 4, 6\}$ y $B = \{1, 2, 3, 4\}$. Entonces $A - B = \{6\}$ y $B - A = \{1, 3\}$.

Dados los sucesos $A, B, C \in \mathcal{P}(E)$, se cumplen las siguientes propiedades:

1. $A \cup A = A$, $A \cap A = A$ (idempotencia).
2. $A \cup B = B \cup A$, $A \cap B = B \cap A$ (conmutativa).
3. $(A \cup B) \cup C = A \cup (B \cup C)$, $(A \cap B) \cap C = A \cap (B \cap C)$ (asociativa).
4. $(A \cup B) \cap C = (A \cap C) \cup (B \cap C)$, $(A \cap B) \cup C = (A \cup C) \cap (B \cup C)$ (distributiva).
5. $A \cup \emptyset = A$, $A \cap E = A$ (elemento neutro).
6. $A \cup E = E$, $A \cap \emptyset = \emptyset$ (elemento absorbente).
7. $A \cup \bar{A} = E$, $A \cap \bar{A} = \emptyset$ (elemento simétrico complementario).
8. $\bar{\bar{A}} = A$ (doble contrario).
9. $\overline{A \cup B} = \bar{A} \cap \bar{B}$, $\overline{A \cap B} = \bar{A} \cup \bar{B}$ (leyes de Morgan).
10. $A \cap B \subseteq A \cup B$.

1.3 Definición de probabilidad

En todo experimento aleatorio existe incertidumbre sobre el resultado de la realización del experimento. La probabilidad trata de cuantificar el grado de incertidumbre asociada a cada suceso de un experimento aleatorio. A lo largo de la historia se han utilizado distintas definiciones del concepto de probabilidad. A continuación se presentan las más comunes.

Definición 6.9 — Probabilidad clásica de Laplace. Para un experimento aleatorio donde todos los elementos del espacio muestral E son equiprobables, se define la *probabilidad* de un suceso $A \subseteq E$ como el cociente entre el número de elementos de A y el número de elementos de E :

$$P(A) = \frac{|A|}{|E|} = \frac{\text{nº casos favorables a } A}{\text{nº casos posibles}}$$

■ **Ejemplo 6.6** Si se considera el espacio muestral correspondiente al lanzamiento de un dado $E = \{1, 2, 3, 4, 5, 6\}$, y el suceso correspondiente a sacar un número par $A = \{2, 4, 6\}$, según la regla de Laplace, la probabilidad de sacar par al tirar un dado es

$$P(A) = \frac{|A|}{|E|} = \frac{3}{6} = 0,5,$$

es decir, un 50 %.

Esta definición es ampliamente utilizada, aunque tiene importantes restricciones:

- No puede utilizarse con espacios muestrales infinitos, o de los que no se conoce el número de casos posibles.
- Es necesario que todos los elementos del espacio muestral tengan la misma probabilidad de ocurrir (*equiprobabilidad*).

Estas restricciones suelen cumplirse en los experimentos relacionados con los juegos de azar (lanzamiento de dados, monedas, etc.) pero es raro que ocurran en los experimentos de las ciencias de la salud. Por ejemplo, los grupos sanguíneos de una población humana no suelen ser equiprobables (normalmente el grupo A es mucho más probable que el resto.)

En estos casos, y gracias al siguiente teorema, es mejor definir la probabilidad a partir de la frecuencia de cada suceso.

Teorema 6.1 — Ley de los grandes números. Cuando un experimento aleatorio se repite un gran número de veces, las frecuencias relativas de los sucesos del experimento tienden a estabilizarse en torno a cierto número, que es precisamente su probabilidad.

Un ejemplo que demuestra el cumplimiento de esta ley puede realizarse tirando múltiples veces una moneda equilibrada y anotando la frecuencia relativa de caras. A medida que se tire más veces la moneda se verá que la frecuencia relativa de caras se va estabilizando en torno a 0,5 que es la probabilidad de sacar cara.

De acuerdo al teorema anterior, se puede dar la siguiente definición

Definición 6.10 — Probabilidad frecuentista. Para un experimento aleatorio reproducible, se define la *probabilidad* de un suceso $A \subseteq E$ como la frecuencia relativa del suceso A en infinitas repeticiones del experimento:

$$P(A) = \lim_{n \rightarrow \infty} \frac{n_A}{n}$$

■ **Ejemplo 6.7** Si en una determinada población existe un 40 % de personas con grupo sanguíneo A, un 30 % de personas con grupo B, un 20 % con grupo O y un 10 % con grupo AB, de acuerdo a la definición frecuentista de probabilidad podemos decir que $P(A) = 0,4$, $P(B) = 0,3$, $P(O) = 0,2$ y $P(AB) = 0,1$.

Aunque esta definición es muy útil en experimentos científicos reproducibles, también tiene serios inconvenientes, ya que

- Sólo se calcula una aproximación de la probabilidad real.
- La repetición del experimento debe ser en las mismas condiciones.

No fue hasta el siglo XX cuando Kolmogórov dió una definición de probabilidad que unificaba las anteriores y que actualmente es la más aceptada.

Definición 6.11 — Kolmogórov. Se llama *probabilidad* a toda aplicación que asocia a cada suceso A del espacio de sucesos de un experimento aleatorio, un número real $P(A)$, que cumple los siguientes axiomas:

1. La probabilidad de un suceso cualquiera es positiva o nula:

$$P(A) \geq 0.$$

2. La probabilidad de la unión de dos sucesos incompatibles es igual a la suma de las probabilidades de cada uno de ellos:

$$P(A \cup B) = P(A) + P(B).$$

3. La probabilidad del suceso seguro es igual a la unidad:

$$P(E) = 1.$$

A partir de los axiomas de la definición de probabilidad se pueden deducir los siguientes resultados:

1. $P(\bar{A}) = 1 - P(A)$.

2. $P(\emptyset) = 0$.
3. Si $A \subseteq B$ entonces $P(A) \leq P(B)$.
4. $P(A) \leq 1$.
5. Si A y B son sucesos compatibles, es decir, su intersección no es vacía, entonces

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

6. Si el suceso A está compuesto por los sucesos elementales e_1, e_2, \dots, e_n , entonces

$$P(A) = \sum_{i=1}^n P(e_i).$$

Este último resultado es especialmente interesante, pues permite calcular la probabilidad de cualquier suceso de manera muy sencilla simplemente sumando las probabilidades de los elementos que lo componen.

1.4 Probabilidad condicionada

La incertidumbre sobre un suceso depende de la información que se tenga sobre el experimento aleatorio. En algunas ocasiones puede que haya que calcular la probabilidad de algún suceso A sabiendo que ha ocurrido otro B . En tal caso se dice que el suceso B es un *condicionante*, y la probabilidad del suceso condicionado suele escribirse como

$$P(A/B)$$

Los condicionantes, en el fondo, cambian el espacio muestral del experimento y por tanto las probabilidades de sus sucesos.

■ **Ejemplo 6.8** Supongamos que hemos observado las siguientes frecuencias de fumadores en un grupo de 100 hombres y 100 mujeres:

	Fuma	No Fuma
Mujeres	30	70
Hombres	40	60
Total	70	130

Entonces, utilizando la definición de frecuentista, la probabilidad de que una persona elegida al azar sea fumadora es la frecuencia relativa de fumar que es $P(\text{Fumar}) = 70/200 = 0,35$.

Sin embargo, si se añade información sobre el experimento y nos dicen que la persona elegida es mujer, entonces la muestra se restringiría sólo a las mujeres y la frecuencia relativa de fumar en mujeres es $P(\text{Fumar}/\text{Mujer}) = 30/100 = 0,3$.

El problema de los condicionamientos es que suelen cambiar el espacio muestral de partida. Afortunadamente, es posible calcular probabilidades condicionadas sin cambiar de espacio muestral gracias a la siguiente fórmula.

Definición 6.12 — Probabilidad condicionada. Dados dos sucesos A y B de un mismo espacio de sucesos de un experimento aleatorio, la probabilidad de A *condicionada* por B es

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

siempre y cuando, $P(B) \neq 0$.

Esta definición permite calcular probabilidades sin tener que alterar el espacio muestral original del experimento.

■ **Ejemplo 6.9** Siguiendo con el anterior, si se calcula la probabilidad de fumar en el caso de ser mujer con esta fórmula se obtiene el mismo resultado:

$$P(\text{Fumar}/\text{Mujer}) = \frac{P(\text{Fumar} \cap \text{Mujer})}{P(\text{Mujer})} = \frac{30/200}{100/200} = \frac{30}{100} = 0,3.$$

De esta definición se deduce que la probabilidad de la intersección es

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

En ocasiones, saber que un determinado suceso ha ocurrido no cambia la incertidumbre sobre otro suceso del mismo experimento. Por ejemplo, si se tiran dos monedas, está claro que el resultado de la primera no cambia la incertidumbre sobre el resultado de la segunda. En tal caso se dice que los sucesos son independientes.

Definición 6.13 — Sucesos independientes. Dados dos sucesos A y B de un mismo espacio de sucesos de un experimento aleatorio, se dice que A es *independiente* de B , si la probabilidad de A no se ve alterada al condicionar por B , es decir,

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

Si A es independiente de B , también se cumple que B es independiente de A , y en general simplemente se dice que A y B son independientes.

También se cumple que si A y B son independientes, entonces

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

1.5 Espacios probabilísticos

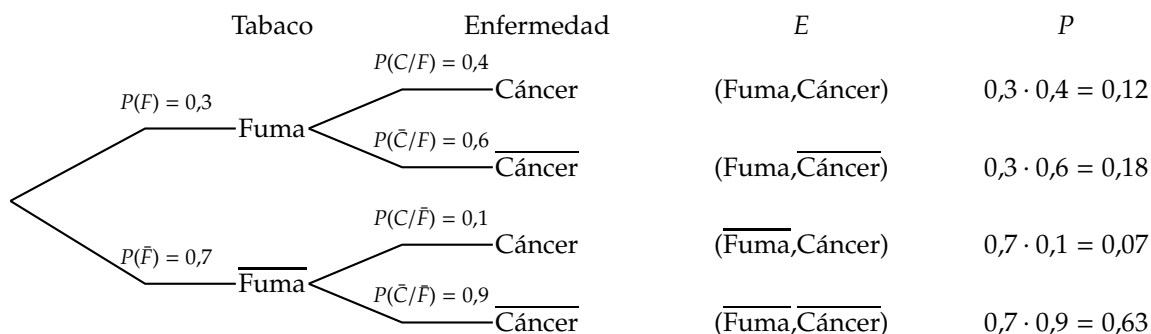
Si a cada uno de los elementos del espacio muestral de un experimento se le asocia su probabilidad, se obtiene un *espacio probabilístico*.

Resulta sencillo construir espacios probabilísticos a partir de los diagramas de árbol que se vieron para construir espacios muestrales. Para ello se deben etiquetar las ramas del árbol con probabilidades del siguiente modo:

1. Para cada nodo del árbol, etiquetar su rama con la probabilidad de que la variable correspondiente tome el valor del nodo, condicionada por la ocurrencia de todos los nodos que conducen hasta el actual.
2. La probabilidad de cada suceso elemental se calcula multiplicando las probabilidades que etiquetan las ramas que conducen hasta él.

■ **Ejemplo 6.10 — Espacio probabilístico con variables dependientes.** Sea una población en la que el 30 % de las personas fuman, y que la incidencia del cáncer de pulmón en fumadores es del 40 % mientras que en los no fumadores es del 10 %.

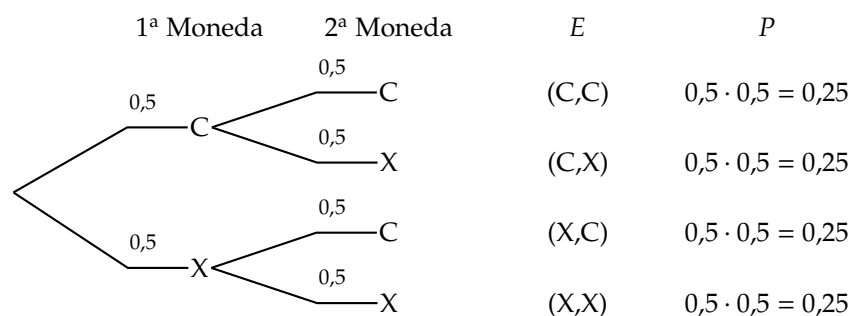
El espacio probabilístico de este experimento es:



Obsérvese que el fumar o no depende del sexo, así que las ramas que salen del suceso mujer no tienen las mismas probabilidades que las que salen del suceso hombre.

Cuando las variables observadas en el experimento son independientes, la construcción del espacio probabilístico se simplifica, ya que las probabilidades condicionadas se convierten en probabilidades simples.

■ **Ejemplo 6.11 — Espacio probabilístico con variables independientes.** El espacio probabilístico del experimento aleatorio que consiste en el lanzamiento de dos monedas es:



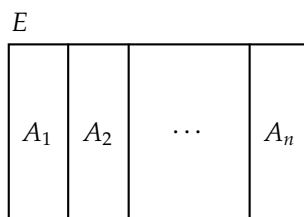
Obsérvese ahora que el resultado de la segunda moneda no depende del resultado de la primera, de manera que las ramas que salen del suceso cara en la primera moneda tienen las mismas probabilidades que las que salen del suceso cruz.

1.6 Teorema de la probabilidad total

En algunos experimentos es posible descomponer el espacio muestral en partes que forman un sistema completo de sucesos.

Definición 6.14 — Sistema completo de sucesos. Una colección de sucesos A_1, A_2, \dots, A_n de un mismo espacio de sucesos es un *sistema completo* si cumple las siguientes condiciones:

1. La unión de todos es el espacio muestral: $A_1 \cup \dots \cup A_n = E$.
2. Son incompatibles dos a dos: $A_i \cap A_j = \emptyset \forall i \neq j$.

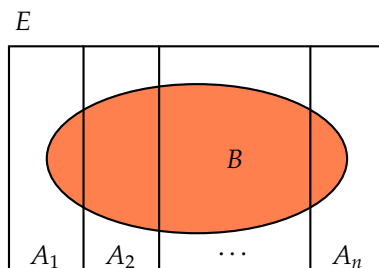


En realidad un sistema completo de sucesos es una partición del espacio muestral de acuerdo a algún atributo, como por ejemplo el sexo o el grupo sanguíneo.

Conocer las probabilidades de un determinado suceso en cada una de las partes de un sistema completo puede ser útil para calcular su probabilidad a partir del siguiente teorema.

Teorema 6.2 — Probabilidad total. Dado un sistema completo de sucesos A_1, \dots, A_n y un suceso B de un mismo espacio de sucesos, se cumple

$$P(B) = \sum_{i=1}^n P(A_i)P(B/A_i).$$



■ **Ejemplo 6.12** Un determinado síntoma B puede ser originado por una enfermedad A pero también lo pueden presentar las personas sin la enfermedad. Se sabe que en la población la tasa de personas con la enfermedad A es $0,2$. Además, de las personas que presentan la enfermedad, el 90% presentan el síntoma, mientras que de las personas sin la enfermedad sólo lo presentan el 40% .

Si se toma una persona al azar de la población, ¿qué probabilidad hay de que tenga el síntoma?

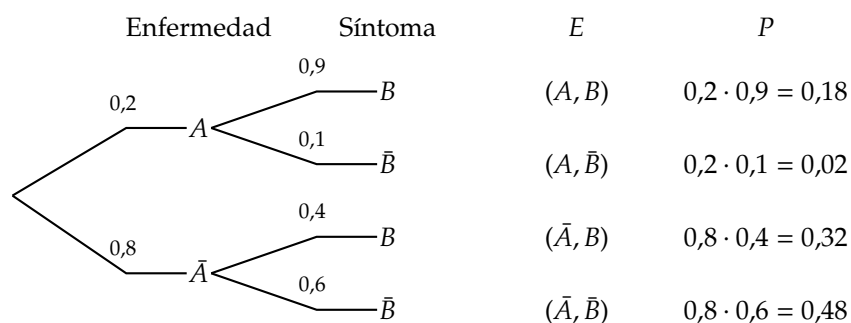
Para responder a la pregunta hay que fijarse en que el conjunto de sucesos $\{A, \bar{A}\}$ es un sistema completo, ya que $A \cup \bar{A} = E$ y $A \cap \bar{A} = \emptyset$, de modo que se puede aplicar el teorema de la probabilidad total se tiene

$$P(B) = P(A)P(B/A) + P(\bar{A})P(B/\bar{A}) = 0,2 \cdot 0,9 + 0,8 \cdot 0,4 = 0,5.$$

Es decir, la mitad de la población tendrá el síntoma.

En el fondo se trata de una media ponderada de probabilidades.

La respuesta a la pregunta anterior es evidente a la luz del espacio probabilístico del experimento.



$$P(B) = P(A, B) + P(\bar{A}, B) = P(A)P(B/A) + P(\bar{A})P(B/\bar{A}) = 0,2 \cdot 0,9 + 0,8 \cdot 0,4 = 0,18 + 0,32 = 0,5.$$

1.7 Teorema de Bayes

Los sucesos de un sistema completo de sucesos A_1, \dots, A_n también pueden verse como las distintas hipótesis ante un determinado hecho B .

En estas condiciones puede ser útil calcular las probabilidades a posteriori $P(A_i/B)$ de cada una de las hipótesis, es decir, una vez se haya cumplido el suceso B . Para ello se utiliza el teorema de Bayes.

Teorema 6.3 — Bayes. Dado un sistema completo de sucesos A_1, \dots, A_n y un suceso B de un mismo espacio de sucesos, se cumple

$$P(A_i/B) = \frac{P(A_i \cap B)}{P(B)} = \frac{P(A_i)P(B/A_i)}{\sum_{i=1}^n P(A_i)P(B/A_i)}.$$

■ **Ejemplo 6.13 — Diagnóstico de una enfermedad.** En el ejemplo anterior se ha visto cómo calcular la probabilidad de que una persona elegida al azar presente el síntoma, pero desde un punto de vista de diagnóstico clínico es mucho más interesante calcular la probabilidad de que una persona que presenta el síntoma tenga la enfermedad, ya que esto permitiría diagnosticar la enfermedad o descartarla.

En este caso, las hipótesis ante las que hay que decidir son A y \bar{A} y sus probabilidades “a priori” son $P(A) = 0,2$ y $P(\bar{A}) = 0,8$.

Esto quiere decir que si no hubiese ninguna información sobre la persona, el diagnóstico sería que no tiene la enfermedad pues es mucho más probable que la tenga.

Sin embargo, si al reconocer a la persona se observa que presenta el síntoma, dicha información condiciona a las hipótesis, y para decidir entre ellas es necesario calcular sus probabilidades "a posteriori", es decir

$$P(A/B) \text{ y } P(\bar{A}/B)$$

Para calcular estas probabilidades se puede utilizar el teorema de Bayes.

$$P(A/B) = \frac{P(A)P(B/A)}{P(A)P(B/A) + P(\bar{A})P(B/\bar{A})} = \frac{0,2 \cdot 0,9}{0,2 \cdot 0,9 + 0,8 \cdot 0,4} = \frac{0,18}{0,5} = 0,36,$$

$$P(\bar{A}/B) = \frac{P(\bar{A})P(B/\bar{A})}{P(A)P(B/A) + P(\bar{A})P(B/\bar{A})} = \frac{0,8 \cdot 0,4}{0,2 \cdot 0,9 + 0,8 \cdot 0,4} = \frac{0,32}{0,5} = 0,64.$$

Según esto, a pesar de que la probabilidad de estar enfermo ha aumentado, seguiríamos diagnosticando que no lo está, puesto que es más probable.

En este caso se dice que el síntoma B no es determinante a la hora de diagnosticar la enfermedad, pues la información que aporta no sirve para cambiar el diagnóstico en ningún caso.

1.8 Tests diagnósticos

En epidemiología es común el uso de tests para diagnosticar enfermedades.

Generalmente estos tests no son totalmente fiables, sino que hay cierta probabilidad de acierto o fallo en el diagnóstico, que suele representarse en la siguiente tabla:

	Presencia de la enfermedad (E)	Ausencia de la enfermedad (\bar{E})
Test positivo (+)	Diagnóstico acertado $P(+/E)$ Sensibilidad	Diagnóstico erróneo $P(+/\bar{E})$
Test negativo (-)	Diagnóstico erróneo $P(-/E)$	Diagnóstico acertado $P(-/\bar{E})$ Especificidad

La validez de una prueba diagnóstica depende de estas dos probabilidades:

Sensibilidad Es el porcentaje de positivos entre las personas enfermas: $P(+/E)$.

Especificidad Es el porcentaje de negativos entre las personas sanas: $P(-/\bar{E})$.

Pero lo realmente interesante de un test diagnóstico es su capacidad predictiva para diagnosticar, lo cual se mide mediante las siguientes probabilidades a posteriori:

Valor predictivo positivo Es el porcentaje de enfermos entre los positivos: $P(E/+)$.

Valor predictivo negativo Es el porcentaje de sanos entre los negativos: $P(\bar{E}/-)$.

Sin embargo, estos últimos valores dependen del porcentaje de enfermos en la población $P(E)$, lo que se conoce como, *tasa o prevalencia* de la enfermedad.

■ **Ejemplo 6.14** Un test para diagnosticar la gripe tiene una sensibilidad del 95 % y una especificidad del 90 %. Según esto, las probabilidades de acierto y fallo del test son:

	Gripe	No gripe
Test +	0,95	0,10
Test -	0,05	0,90

Si la prevalencia de la gripe en la población es del 10 % y al aplicar el test a un individuo da positivo, ¿cuál es la probabilidad de que tenga gripe?

Aplicando el teorema de Bayes, se tiene que el valor predictivo positivo del test vale

$$P(\text{Gripe}/+) = \frac{P(\text{Gripe})P(+/\text{Gripe})}{P(\text{Gripe})P(+/\text{Gripe}) + P(\bar{\text{Gripe}})P(+/\bar{\text{Gripe}})} = \frac{0,1 \cdot 0,95}{0,1 \cdot 0,95 + 0,9 \cdot 0,1} = 0,5135.$$

La probabilidad de no tener la gripe sería $P(\overline{\text{Gripe}}) = 1 - 0,51 = 0,49$, y como la probabilidad de tener la gripe es mayor que la de no tenerla, se diagnosticaría que tiene gripe. Sin embargo, el valor predictivo positivo de este test es muy bajo y nos equivocaríamos en el 49 % de los diagnósticos de enfermedad.

Y si el test da negativo, ¿cuál es la probabilidad de que no tenga gripe?

De nuevo, aplicando el teorema de Bayes, se tiene que el valor predictivo negativo del test vale

$$P(\overline{\text{Gripe}}/-) = \frac{P(\overline{\text{Gripe}})P(-/\overline{\text{Gripe}})}{P(\text{Gripe})P(-/\text{Gripe}) + P(\overline{\text{Gripe}})P(-/\overline{\text{Gripe}})} = \frac{0,9 \cdot 0,9}{0,1 \cdot 0,05 + 0,9 \cdot 0,9} = 0,9939.$$

De manera que el valor predictivo negativo de este test es mucho más alto que el valor predictivo positivo, y por tanto este test es mucho más útil para descartar la enfermedad que para detectarla.

2 Ejercicios resueltos

1. Construir los espacios muestrales correspondientes a los siguientes experimentos aleatorios:

a) Sacar una carta de una baraja española.



- 1) Seleccionar el menú Teaching ▶ Probabilidad ▶ Juegos de azar ▶ Naipes ▶ Espacio probabilístico.
- 2) En el cuadro de diálogo que aparece, desactivar la opción Incluir probabilidades y hacer clic en el botón Enviar.

b) Lanzar dos monedas.



- 1) Seleccionar el menú Teaching ▶ Probabilidad ▶ Juegos de azar ▶ Monedas ▶ Espacio probabilístico.
- 2) En el cuadro de diálogo que aparece, introducir 2 en el campo Número de monedas, desactivar la opción Incluir probabilidades y hacer clic en el botón Enviar.

c) Lanzar dos dados.



- 1) Seleccionar el menú Teaching ▶ Probabilidad ▶ Juegos de azar ▶ Dados ▶ Espacio probabilístico.
- 2) En el cuadro de diálogo que aparece, introducir 2 en el campo Número de dados, desactivar la opción Incluir probabilidades y hacer clic en el botón Enviar.

d) Lanzar dos dados y dos monedas.



- 1) Seleccionar el menú Teaching ▶ Probabilidad ▶ Combinar espacios probabilísticos independientes.
- 2) En el cuadro de diálogo que aparece, seleccionar los conjuntos de datos correspondientes a los espacios muestrales del lanzamiento de dos monedas y del lanzamiento de dos dados generados en los apartados anteriores, y hacer clic en el botón Enviar.

2. Repetir el experimento de lanzar dos monedas 10 veces, 100 veces, 1000 veces y 1000000 de veces y calcular las frecuencias relativas de cada resultado. ¿Hacia dónde tienden las frecuencias? Construir el espacio probabilístico de este experimento y comprobar que se cumple la ley de los grandes números, es decir, que las frecuencias anteriores se aproximan a las probabilidades de cada suceso elemental.



Para la realización del experimento:

- a) Seleccionar el menú Teaching ▶ Probabilidad ▶ Juegos de azar ▶ Monedas ▶ Lanzamiento de monedas.
- b) En el cuadro de diálogo que aparece, introducir 2 en el campo Número de monedas, introducir 10 en el campo Número de lanzamientos, activar la casilla Distribución de frecuencias y hacer clic en el botón Enviar.

Repetir los mismos pasos pero introduciendo 100, 1000 y 1000000 respectivamente en el campo Número de lanzamientos.

Para construir el espacio probabilístico correspondiente:

- a) Seleccionar el menú Teaching ▶ Probabilidad ▶ Juegos de azar ▶ Monedas ▶ Espacio probabilístico.

- b) En el cuadro de diálogo que aparece, introducir 2 en el campo Número de monedas y hacer clic en el botón Enviar.

3. En una estantería hay tres cajas de un medicamento A, dos de un medicamento B y una de un medicamento C. Construir los espacios probabilísticos asociados a los siguientes experimentos aleatorios:

- a) Elegir tres medicamentos al azar sin reemplazamiento.



- 1) Seleccionar el menú Teaching ▶ Probabilidad ▶ Juegos de azar ▶ Urna ▶ Espacio probabilístico.
- 2) En el cuadro de diálogo que aparece, seleccionar la opción Lista de objetos, introducir los objetos A,A,A,B,B,C en el campo Lista de objetos, introducir 3 en el campo Número de extracciones, y hacer clic en el botón Enviar.

- b) Elegir tres medicamentos al azar con reemplazamiento.



Repetir los mismos pasos del apartado anterior pero además activando la casilla Con reemplazamiento.

4. Gregor Mendel, monje austríaco, desarrollo en el siglo XIX los principios fundamentales de genética. Mendel demostró que las características heredables se transmiten en unidades discretas que se heredan por separado en cada generación. Estas unidades discretas, que Mendel llamó *elemente*, se conocen hoy como *genes*.

Cada característica hereditaria depende de dos factores separados que provienen uno de cada progenitor. Estos factores son los *alelos* de cada gen, que pueden ser *dominantes* (cuando se expresan en el fenotipo sin tener en cuenta el otro alelo) o *recesivos* (que se expresa sólo cuando el otro alelo es igual).

Mendel demostró que en la reproducción los alelos se combinan aleatoriamente y de manera independiente para formar el gen del hijo.

En uno de sus experimentos cruzó dos plantas de guisantes con idéntico genotipo Aa y Bb, donde el primer gen se refiere al color del guisante (A amarillo y a verde) y el segundo gen se refiere a la forma del guisante (B liso y b rugoso). Se pide:

- a) Construir el espacio probabilístico correspondiente al genotipo del gen del color. ¿Cuál es la probabilidad que la planta resultante diese guisantes con fenotipo amarillo? ¿Y verde?



Para construir el espacio probabilístico:

- 1) Crear un conjunto de datos mendel.color con la variable alelo.color.progenitor con los dos posibles alelos del progenitor correspondientes al gen del color (A y a).
- 2) Seleccionar el menú Teaching ▶ Probabilidad ▶ Construir espacio probabilístico.
- 3) En el cuadro de diálogo que aparece seleccionar el conjunto de datos mendel.color, darle el nombre mendel.color.ep al espacio probabilístico y hacer clic en el botón Enviar.
- 4) Seleccionar el menú Teaching ▶ Probabilidad ▶ Repetir espacio probabilístico.
- 5) En el cuadro de diálogo que aparece seleccionar el conjunto de datos mendel.color.ep, darle el mismo nombre al espacio probabilístico resultante y hacer clic en el botón Enviar.

Para calcular la probabilidad de fenotipo amarillo:

- 1) Seleccionar el menú Teaching ▶ Probabilidad ▶ Calcular probabilidad.

- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico `mendel.color.ep`, introducir `alelo.color.progenitor1 == "A" | alelo.color.progenitor2 == "A"` en el campo Suceso y hacer clic en el botón Enviar.

Para calcular la probabilidad de fenotipo verde:

- 1) Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico `mendel.color.ep`, introducir `alelo.color.progenitor1 == "a" & alelo.color.progenitor2 == "a"` en el campo Suceso y hacer clic en el botón Enviar.

- b) Construir el espacio probabilístico correspondiente al genotipo de ambos genes. ¿Cuál es la probabilidad que la planta resultante diese guisantes con fenotipo amarillo y liso? ¿Y verde y rugoso?



Para construir el espacio probabilístico:

- 1) Crear un conjunto de datos `mendel.forma` con la variable `alelo.forma.progenitor` con los dos posibles alelos del progenitor correspondientes al gen de la forma (B y b).
- 2) Seleccionar el menú Teaching►Probabilidad►Construir espacio probabilístico.
- 3) En el cuadro de diálogo que aparece seleccionar el conjunto de datos `mendel.forma`, darle el nombre `mendel.forma.ep` al espacio probabilístico y hacer clic en el botón Enviar.
- 4) Seleccionar el menú Teaching►Probabilidad►Repetir espacio probabilístico.
- 5) En el cuadro de diálogo que aparece seleccionar el conjunto de datos `mendel.forma.ep`, darle el mismo nombre al espacio probabilístico resultante y hacer clic en el botón Enviar.
- 6) Seleccionar el menú Teaching►Probabilidad►Combinar espacios probabilísticos independientes.
- 7) En el cuadro de diálogo que aparece, seleccionar los conjuntos de datos correspondientes a los espacios probabilísticos `mendel.color.ep` y `mendel.forma.ep`, y hacer clic en el botón Enviar.

Para calcular la probabilidad de fenotipo amarillo y liso:

- 1) Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico `mendel.color.ep`, introducir `(alelo.color.progenitor1 == "A" | alelo.color.progenitor2 == "A") & (alelo.forma.progenitor1 == "B" | alelo.forma.progenitor2 == "B")` en el campo Suceso y hacer clic en el botón Enviar.

Para calcular la probabilidad de fenotipo verde y rugoso:

- 1) Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico `mendel.color.ep`, introducir `alelo.color.progenitor1 == "a" & alelo.color.progenitor2 == "a" & alelo.forma.progenitor1 == "a" & alelo.forma.progenitor2 == "a"` en el campo Suceso y hacer clic en el botón Enviar.

5. En una población se ha hecho un estudio epidemiológico sobre tres enfermedades asociadas habitualmente a la infancia, como son la varicela, el sarampión y la rubeola. Las frecuencias observadas aparecen en la siguiente tabla:

Varicela	Sarampión	Rubeola	Frecuencia
No	No	No	2654
No	No	Si	1436
No	Si	No	1682
No	Si	Si	668
Si	No	No	1747
Si	No	Si	476
Si	Si	No	876
Si	Si	Si	265

Se pide:

- Crear el conjunto de datos enfermedades.infantiles con las variables varicela, sarampion, rubeola y frecuencia e introducir datos de la población.
- Crear el espacio probabilístico asociado a la población.



- Seleccionar el menú Teaching►Probabilidad►Construir espacio probabilístico.
- En el cuadro de diálogo que aparece seleccionar el conjunto de datos enfermedades.infantiles, activar la casilla Definir frecuencias, seleccionar la variable frecuencia en el campo Frecuencia, darle el nombre enfermedades.infantiles.ep al espacio probabilístico y hacer clic en el botón Enviar.

- Calcular la probabilidad de que una persona de esta población haya tenido varicela.



- Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- En el cuadro de diálogo que aparece seleccionar el espacio probabilístico enfermedades.infantiles.ep, introducir varicela == "Si" en el campo Suceso y hacer clic en el botón Enviar.

- Calcular la probabilidad de que una persona de esta población haya tenido varicela o sarampión.



- Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- En el cuadro de diálogo que aparece seleccionar el espacio probabilístico enfermedades.infantiles.ep, introducir varicela == "Si" | sarampion=="Si" en el campo Suceso y hacer clic en el botón Enviar.

- Calcular la probabilidad de que una persona de esta población haya tenido sarampión y rubeola.



- Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- En el cuadro de diálogo que aparece seleccionar el espacio probabilístico enfermedades.infantiles.ep, introducir sarampion == "Si" & rubeola=="Si" en el campo Suceso y hacer clic en el botón Enviar.

- Calcular la probabilidad de que una persona de esta población haya tenido varicela si no ha tenido sarampion. ¿Son independientes el haber tenido varicela y el haber tenido sarampión?



- Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- En el cuadro de diálogo que aparece seleccionar el espacio probabilístico enfermedades.infantiles.ep, introducir varicela == "Si" en el campo Suceso, activar la casilla

Probabilidad condicionada, introducir `sarampion == "No"` en el campo Condición y hacer clic en el botón Enviar.

- g) Calcular la probabilidad de que una persona de esta población no haya tenido rubeola ni sarampión si ha tenido varicela.



- 1) Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico `enfermedades.infantiles.ep`, introducir `rubeola == "No"` & `sarampion=="No"` en el campo Suceso, activar la casilla Probabilidad condicionada, introducir `varicela == "Si"` en el campo Condición y hacer clic en el botón Enviar.

6. Se ha probado un test diagnóstico para detectar el embarazo en un grupo de mujeres en edad de procrear, obteniendo los siguientes resultados

Embarazo	Test	Frecuencia
No	–	3876
No	+	47
Si	–	12
Si	+	131

Se pide:

- a) Crear el conjunto de datos `test.embarazo` con las variables `embarazo`, `test`, y `frecuencia` e introducir datos de la muestra.
- b) Crear el espacio probabilístico asociado a la población.



- 1) Seleccionar el menú Teaching►Probabilidad►Construir espacio probabilístico.
- 2) En el cuadro de diálogo que aparece seleccionar el conjunto de datos `test.embarazo`, activar la casilla Definir frecuencias, seleccionar la variable frecuencia en el campo Frecuencia, darle el nombre `test.embarazo.ep` al espacio probabilístico y hacer clic en el botón Enviar.

- c) Calcular la prevalencia del embarazo.



- 1) Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico `test.embarazo.ep`, introducir `embarazo == "Si"` en el campo Suceso y hacer clic en el botón Enviar.

- d) Calcular la probabilidad de que el test de positivo.



- 1) Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico `test.embarazo.ep`, introducir `test == "+"` en el campo Suceso y hacer clic en el botón Enviar.

- e) Calcular la sensibilidad del test.



- 1) Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico `test.embarazo.ep`, introducir `test=="+"` en el campo Suceso, activar la casilla Probabilidad condicionada, introducir `embarazo == "Si"` en el campo Condición y hacer clic en el botón Enviar.

- f) Calcular la especificidad del test.



- 1) Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico test.embarazo.ep, introducir test=="-" en el campo Suceso, activar la casilla Probabilidad condicionada, introducir embarazo == "No" en el campo Condición y hacer clic en el botón Enviar.

g) Calcular el valor predictivo positivo del test. ¿Es útil el test para detectar la enfermedad?



- 1) Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico test.embarazo.ep, introducir embarazo=="Si" en el campo Suceso, activar la casilla Probabilidad condicionada, introducir test=="+" en el campo Condición y hacer clic en el botón Enviar.

h) Calcular el valor predictivo negativo del test. ¿Es útil el test para descartar la enfermedad?



- 1) Seleccionar el menú Teaching►Probabilidad►Calcular probabilidad.
- 2) En el cuadro de diálogo que aparece seleccionar el espacio probabilístico test.embarazo.ep, introducir embarazo=="No" en el campo Suceso, activar la casilla Probabilidad condicionada, introducir test=="-" en el campo Condición y hacer clic en el botón Enviar.

3 Ejercicios propuestos

1. Construir el espacio muestral correspondiente a tirar una moneda, un dado y sacar una carta de una baraja española.
2. Para comprobar la eficacia de una vacuna contra la gripe se tomó una muestra de 1000 y se observó si fueron vacunadas y si finalmente tuvieron gripe o no. Los resultados obtenidos fueron los siguientes

Vacuna	Gripe	Frecuencia
No	No	418
No	Si	312
Si	No	233
Si	Si	37

Se pide:

- a) Construir el espacio probabilístico asociado al experimento.
 - b) Calcular la probabilidad de haberse vacunado contra la gripe.
 - c) Calcular la prevalencia de la gripe.
 - d) Calcular la probabilidad de desarrollar la gripe tras haberse vacunado. ¿Es efectiva la vacuna?
3. Para probar la eficacia de un test diagnóstico para detectar el ébola en un país centroafricano, se tomó una muestra de personas a las que se le ha aplicado el test. El test dió positivo en 147 personas con ébola, pero también dió positivo en 28 personas sin ébola. Por otro lado el test dió negativo en 97465 personas sin ébola, pero también dió negativo en 65 personas con ébola. Se pide:
- a) Construir el espacio probabilístico asociado al test diagnóstico.

- b) Calcula la prevalencia del ébola en ese país.
- c) Calcular la probabilidad de que el test de negativo.
- d) Calcular la sensibilidad y la especificidad del test.
- e) ¿Para qué es más efectivo el test, para detectar o para descartar el ébola?

Variables Aleatorias Discretas

1 Fundamentos teóricos

1.1 Variables Aleatorias

Se define una *variable aleatoria* asignando a cada resultado del experimento aleatorio un número. Esta asignación puede realizarse de distintas maneras, obteniéndose de esta forma diferentes variables aleatorias. Así, en el lanzamiento de dos monedas podemos considerar el número de caras o el número de cruces. En general, si los resultados del experimento son numéricos, se tomarán dichos números como los valores de la variable, y si los resultados son cualitativos, se hará corresponder a cada modalidad un número arbitrariamente.

Formalmente, una *variable aleatoria* X es una función real definida sobre los puntos del espacio muestral E de un experimento aleatorio.

$$X : E \rightarrow \mathbb{R}$$

De esta manera, la distribución de probabilidad del espacio muestral E , se transforma en una distribución de probabilidad para los valores de X .

El conjunto formado por todos los valores distintos que puede tomar la variable aleatoria se llama *Rango* o *Recorrido* de la misma.

Las variables aleatorias pueden ser de dos tipos: discretas o continuas. Una variable es *discreta* cuando sólo puede tomar valores aislados, mientras que es *continua* si puede tomar todos los valores posibles de un intervalo.

1.2 Variables Aleatorias Discretas (v.a.d.)

Se considera una v.a.d. X que puede tomar los valores $x_i, i = 1, 2, \dots, n$.

Función de probabilidad

La *distribución de probabilidad* de X se suele caracterizar mediante una función $f(x)$, conocida como *función de probabilidad*, que asigna a cada valor de la variable su probabilidad. Esto es

$$f(x_i) = P(X = x_i), \quad i = 1, \dots, n$$

verificándose que

$$\sum_{i=1}^n f(x_i) = 1$$

Función de distribución

Otra forma equivalente de caracterizar la distribución de probabilidad de X es mediante otra función $F(x)$, llamada *función de distribución*, que asigna a cada $x \in \mathbb{R}$ la probabilidad de que X tome un valor menor o igual que dicho número x . Así,

$$F(x) = P(X \leq x) = \sum_{x_i \leq x} f(x_i)$$

Tanto la función de probabilidad como la función de distribución pueden representarse de forma gráfica, tal y como se muestra en la figura 7.1.

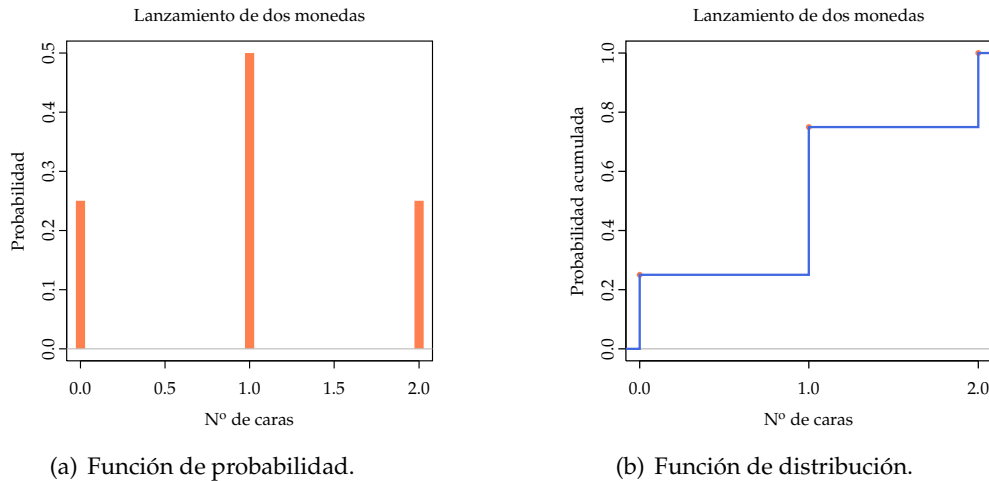


Figura 7.1 – Función de probabilidad y función de distribución de la variable aleatoria X que mide el número de caras obtenido en el lanzamiento de dos monedas.

Estadísticos poblacionales

Los parámetros descriptivos más importantes de una v.a.d. X son:

Media o Esperanza

$$E[X] = \mu = \sum_{i=1}^n x_i f(x_i)$$

Varianza

$$V[X] = \sigma^2 = \sum_{i=1}^n (x_i - \mu)^2 f(x_i) = \sum_{i=1}^n x_i^2 f(x_i) - \mu^2$$

Desviación típica

$$D[X] = \sigma = +\sqrt{\sigma^2}$$

La media es una medida de tendencia central, mientras que la varianza y la desviación típica son medidas de dispersión.

Entre las v.a.d. cabe destacar las denominadas *Binomial* y de *Poisson*.

Variable Binomial

Se considera un experimento aleatorio en el que puede ocurrir el suceso A o su contrario \bar{A} , con probabilidades p y $1 - p$ respectivamente.

Si se realiza el experimento anterior n veces, la v.a.d. X que recoge el número de veces que ha ocurrido el suceso A , se denomina *Variable Binomial* y se designa por $X \sim B(n, p)$.

El recorrido de la variable X es $\{0, 1, \dots, n\}$ y su función de probabilidad viene dada por

$$f(x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

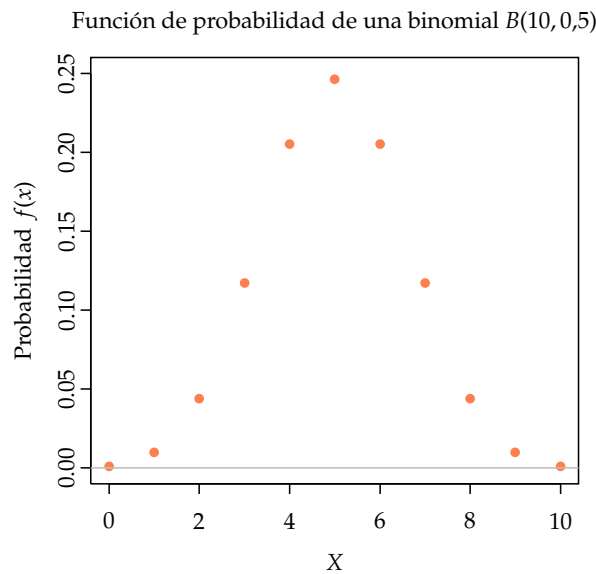


Figura 7.2 – Función de probabilidad de una variable aleatoria binomial de 10 repeticiones y probabilidad de éxito 0.5

cuya gráfica se puede apreciar en la figura 7.2.

A partir de la expresión anterior se puede demostrar que

$$\begin{aligned}\mu &= np \\ \sigma^2 &= np(1-p) \\ \sigma &= +\sqrt{np(1-p)}\end{aligned}$$

En el caso particular de que el experimento se realice una sola vez, la variable aleatoria recibe el nombre de Variable de Bernoulli. Una variable Binomial $X \sim B(n, p)$ se puede considerar como suma de n variables de Bernoulli idénticas con distribución $B(1, p)$.

Variable de Poisson

Las variables de Poisson surgen de la observación de un conjunto discreto de fenómenos puntuales en un soporte continuo de tiempo, longitud o espacio. Por ejemplo: nº de llamadas que llegan a una centralita telefónica en un tiempo establecido, nº de hematíes en un volumen de sangre, etc. Se supone además que en un soporte continuo suficientemente grande, el número medio de fenómenos ocurridos por unidad de soporte considerado, es una constante que designaremos por λ .

A la v.a.d. X , que recoge el número de fenómenos puntuales que ocurren en un intervalo de amplitud establecida, se le denomina *Variable de Poisson* y se designa por $X \sim P(\lambda)$.

El recorrido de la variable X es $\{0, 1, 2, \dots\}$, no existiendo un valor máximo que pueda alcanzar. Su función de probabilidad viene dada por

$$f(x) = \frac{\lambda^x}{x!} e^{-\lambda}$$

y su gráfica aparece en la figura 7.3

Se puede demostrar que

$$\begin{aligned}\mu &= \lambda \\ \sigma^2 &= \lambda \\ \sigma &= +\sqrt{\lambda}\end{aligned}$$

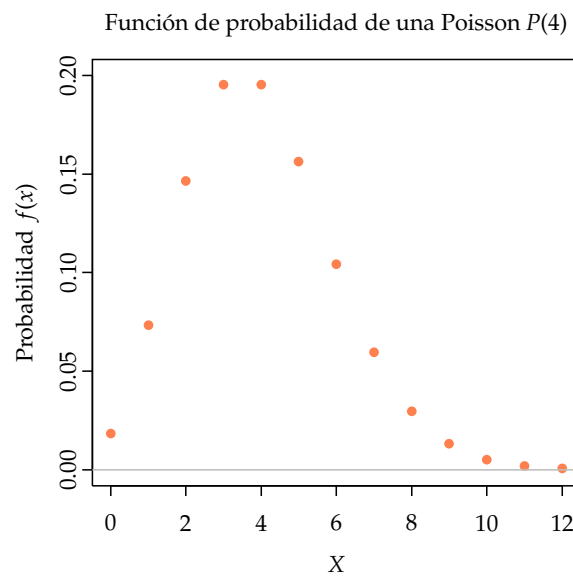


Figura 7.3 – Función de probabilidad de una variable aleatoria Poisson de media $\lambda = 4$

La distribución de Poisson aparece como límite de la distribución Binomial cuando el número n de repeticiones del experimento es muy grande y la probabilidad p de que ocurra el suceso A considerado es muy pequeña. Por ello, la distribución de Poisson se llama también *Ley de los Casos Raros*. En la práctica se considera aceptable realizar los cálculos de probabilidades correspondientes a una variable $B(n, p)$ mediante las fórmulas correspondientes a una variable $P(\lambda)$ con $\lambda = np$, siempre que $n \geq 50$ y $p < 0,1$.

2 Ejercicios resueltos

1. Sea X la variable que mide el número de caras obtenidas al lanzar 10 monedas. Para ver de manera experimental la distribución de probabilidad de X se realiza un experimento aleatorio que consiste en lanzar varias veces las 10 monedas y anotar el número de caras obtenido en cada lanzamiento. Se pide:

- a) Lanzar las 10 monedas 1000 veces y calcular las frecuencias relativas de las caras obtenidas y el diagrama de barras asociado.



Para generar los lanzamientos de monedas:

- 1) Seleccionar el menú Teaching►Simulaciones►Lanzamiento de monedas.
- 2) En el cuadro de diálogo que aparece, introducir 10 en el campo Número de monedas, 1000 en el campo Número de lanzamientos, introducir un nombre para el conjunto de datos y hacer click en el botón aceptar.

Para calcular las frecuencias relativas:

- 1) Seleccionar el menú Teaching►Distribución de frecuencias►Tabla de frecuencias.
- 2) En el cuadro de diálogo que aparece, seleccionar como variable a tabular la variable sum y hacer click en el botón Aceptar.

Para dibujar el diagrama de barras:

- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de barras.
- 2) En el cuadro de diálogo que aparece seleccionar la variable sum.
- 3) En la solapa Opciones de las barras marcar la opción Frecuencias relativas y hacer click en el botón Aceptar.

- b) Generar la distribución de probabilidad de una variable Binomial $B(10, 0,5)$ y compararla con la distribución de frecuencias relativas del apartado anterior.



- 1) Seleccionar el menú teaching►Distribuciones►Discretas►Binomial►Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir 0,1,2,3,4,5,6,7,8,9,10 en el campo Valores de la variable, introducir 10 en el campo Número de repeticiones, 0,5 en el campo Probabilidad de éxito, y hacer click en el botón Aceptar.

- c) Dibujar la gráfica de la función de probabilidad de la Binomial $X \sim B(10, 0,5)$ y compararla con el diagrama de barras de frecuencias relativas del primer apartado.



- 1) Seleccionar el menú Teaching►Distribuciones►Discretas►Binomial►Gráfico de probabilidad.
- 2) En el cuadro de diálogo que aparece, introducir 10 en el campo Número de repeticiones, 0,5 en el campo Probabilidad de éxito y hacer click en el botón Aceptar.

- d) Dibujar la gráfica de la función de distribución.



- 1) Seleccionar el menú Teaching►Distribuciones►Discretas►Binomial►Gráfico de probabilidad.
- 2) En el cuadro de diálogo que aparece, introducir 10 en el campo Número de repeticiones, 0,5 en el campo Probabilidad de éxito, seleccionar la opción Función de distribución y hacer click en el botón Aceptar.

- e) Calcular $P(X = 7)$.



- 1) Seleccionar el menú **teaching** ▶ **Distribuciones** ▶ **Discretas** ▶ **Binomial** ▶ **Probabilidades**.
- 2) En el cuadro de diálogo que aparece, introducir 7 en el campo **Valores de la variable**, introducir 10 en el campo **Número de repeticiones**, 0,5 en el campo **Probabilidad de éxito**, y hacer click en el botón **Aceptar**.

f) Calcular $P(X \leq 4)$.



- 1) Seleccionar el menú **Teaching** ▶ **Distribuciones** ▶ **Discretas** ▶ **Binomial** ▶ **Probabilidades**.
- 2) En el cuadro de diálogo que aparece, introducir 4 en el campo **Valores de la variable**, 10 en el campo **Número de repeticiones**, 0,5 en el campo **Probabilidad de éxito**.
- 3) Seleccionar la opción **Probabilidades acumuladas** y hacer click en el botón **Aceptar**.

g) Calcular $P(X > 5)$.



- 1) Seleccionar el menú **Teaching** ▶ **Distribuciones** ▶ **Discretas** ▶ **Binomial** ▶ **Probabilidades**.
- 2) En el cuadro de diálogo que aparece, introducir 5 en el campo **Valores de la variable**, 10 en el campo **Número de repeticiones**, 0,5 en el campo **Probabilidad de éxito**.
- 3) Seleccionar la opción **Probabilidades acumuladas**, seleccionar la opción **derecha** en el campo **cola de acumulación** y hacer click en el botón **Aceptar**.

h) Calcular $P(2 \leq X < 9)$.



- 1) Seleccionar el menú **Teaching** ▶ **Distribuciones** ▶ **Discretas** ▶ **Binomial** ▶ **Probabilidades**.
 - 2) En el cuadro de diálogo que aparece, introducir los valores 1, 8 en el campo **Valores de la variable**, 10 en el campo **Número de repeticiones**, 0,5 en el campo **Probabilidad de éxito**.
 - 3) Seleccionar la opción **Probabilidades acumuladas** y hacer click en el botón **Aceptar**.
- La probabilidad del intervalo $P(2 \leq X < 9)$ es la resta de las probabilidades obtenidas $P(X < 9) = P(X \leq 8)$ y $P(X < 2) = P(X \leq 1)$.

2. El número de nacimientos diarios en una determinada población sigue una distribución de Poisson de media 6 nacimientos al día. Se pide:

a) Dibujar la gráfica de la función de probabilidad.



- 1) Seleccionar el menú **Teaching** ▶ **Distribuciones** ▶ **Discretas** ▶ **Poisson** ▶ **Gráfico de probabilidad**.
- 2) En el cuadro de diálogo que aparece, introducir el valor 6 en el campo **Media** y hacer click en el botón **Aceptar**.

b) Dibujar la gráfica de la función de distribución.



- 1) Seleccionar el menú **Teaching** ▶ **Distribuciones** ▶ **Discretas** ▶ **Poisson** ▶ **Gráfico de probabilidad**.
- 2) En el cuadro de diálogo que aparece, introducir el valor 6 en el campo **Media**, marcar la opción **Función de distribución** y hacer click en el botón **Aceptar**.

c) Calcular la probabilidad de un día haya 1 nacimiento.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Discretas ▶ Poisson ▶ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir el valor 1 en el campo Valores de la variable, introducir el valor 6 en el campo Media, y hacer click en el botón Aceptar.

d) Calcular la probabilidad de que un día haya menos de 6 nacimientos.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Discretas ▶ Poisson ▶ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir 5 en el campo Valores de la variable y 6 en el campo Media.
- 3) Seleccionar la opción Probabilidades acumuladas y hacer click en el botón Aceptar.

e) Calcular la probabilidad de que un día haya 4 o más nacimientos.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Discretas ▶ Poisson ▶ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir 3 en el campo Valores de la variable y 6 en el campo Media.
- 3) Seleccionar la opción Probabilidades acumuladas, seleccionar la opción derecha en el campo cola de acumulación y hacer click en el botón Aceptar.

f) Calcular la probabilidad de que un día haya entre 4 y 8 nacimientos, inclusivos.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Discretas ▶ Poisson ▶ Probabilidades.
 - 2) En el cuadro de diálogo que aparece, introducir los valores 3, 8 en el campo Valores de la variable y 6 en el campo Media.
 - 3) Seleccionar la opción Probabilidades acumuladas y hacer click en el botón Aceptar.
- La probabilidad del intervalo $P(4 \leq X \leq 8)$ es la resta de las probabilidades obtenidas $P(X \leq 8)$ y $P(X < 4) = P(X \leq 3)$.

3. La ley de los casos raros dice que en una distribución Binomial $B(n, p)$, cuando $n \geq 30$ y $p \leq 0,1$ la distribución se parece mucho a una distribución Poisson $P(np)$. Para comprobar hasta qué punto se parecen estas distribuciones, se pide:

a) Generar la distribución de probabilidad de una variable Binomial $B(30, 0,1)$.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Discretas ▶ Binomial ▶ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir los valores 0,1,2,3,4,5,6,7,8,9,10 en el campo Valores de la variable, introducir el valor 30 en el campo Número de repeticiones, 0,1 en el campo Probabilidad de éxito y hacer click en el botón Aceptar.

b) Generar la distribución de probabilidad de una variable Poisson $P(3)$ y compararla con la de la binomial $B(30, 0,1)$.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Discretas ▶ Poisson ▶ Probabilidades.

- 2) En el cuadro de diálogo que aparece, introducir los valores 0,1,2,3,4,5,6,7,8,9,10 en el campo Valores de la variable, introducir el valor 3 en el campo Media y hacer click en el botón Aceptar.
- c) Generar la distribución de probabilidad de una variable Binomial $B(100, 0,03)$ y compararla con la de la Poisson $P(3)$. ¿Se parecen más estas distribuciones que las anteriores?



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Discretas ▶ Binomial ▶ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir los valores 0,1,2,3,4,5,6,7,8,9,10 en el campo Valores de la variable, introducir el valor 100 en el campo Número de repeticiones, 0,03 en el campo Probabilidad de éxito y hacer click en el botón Aceptar.

- d) Dibujar las gráficas de las distribuciones anteriores y ver cuáles se parecen más. ¿Se cumple la ley de los casos raros?



- 1) Seleccionar el menú Teaching ▶ Simulaciones ▶ Ley de los casos raros.
- 2) En el cuadro de diálogo que aparece, desplazar el deslizador de n hasta 30 y el de p hasta 0.1.
- 3) Después desplazar el deslizador de n hasta 100 y el de p hasta 0.03.

3 Ejercicios propuestos

1. Al lanzar 100 veces una moneda, ¿cuál es la probabilidad de obtener entre 40 y 60 caras inclusive?
2. La probabilidad de curación de un paciente al ser sometido a un determinado tratamiento es 0,85. Calcular la probabilidad de que en un grupo de 6 enfermos sometidos a tratamiento:
 - a) Se curen la mitad.
 - b) Se curen al menos 4.
3. La probabilidad de que al administrar una vacuna dé una determinada reacción es 0,001. Si se vacunan 2000 personas ¿cuál es la probabilidad de que aparezca alguna reacción adversa?
4. El número medio de llamadas por minuto que llegan a una centralita telefónica es igual a 120. Se pide:
 - a) Dar la distribución de probabilidad del número de llamadas en 2 segundos y dibujar su gráfica.
 - b) Calcular la probabilidad de que durante 2 segundos lleguen a la centralita menos de 4 llamadas.
 - c) Calcular la probabilidad de que durante 3 segundos lleguen a la centralita 3 llamadas como mínimo.
5. Se sabe que la probabilidad de que aparezca una bacteria en un mm^3 de cierta disolución es de 0,002. Si en cada mm^3 a lo sumo puede aparecer una bacteria, determinar la probabilidad de que en un cm^3 haya como máximo 5 bacterias.

Variables Aleatorias Continuas

1 Fundamentos teóricos

1.1 Variables Aleatorias

Se define una *variable aleatoria* asignando a cada resultado del experimento aleatorio un número. Esta asignación puede realizarse de distintas maneras, obteniéndose de esta forma diferentes variables aleatorias. Así, en el lanzamiento de dos monedas podemos considerar el número de caras o el número de cruces. En general, si los resultados del experimento son numéricos, se tomarán dichos números como los valores de la variable, y si los resultados son cualitativos, se hará corresponder a cada modalidad un número arbitrariamente.

Formalmente, una *variable aleatoria* X es una función real definida sobre los puntos del espacio muestral E de un experimento aleatorio.

$$X : E \rightarrow \mathbb{R}$$

De esta manera, la distribución de probabilidad del espacio muestral E , se transforma en una distribución de probabilidad para los valores de X .

El conjunto formado por todos los valores distintos que puede tomar la variable aleatoria se llama *Rango* o *Recorrido* de la misma.

Las variables aleatorias pueden ser de dos tipos: discretas o continuas. Una variable es *discreta* cuando sólo puede tomar valores aislados, mientras que es *continua* si puede tomar todos los valores posibles de un intervalo.

1.2 Variables Aleatorias Continuas (v.a.c.)

Se considera una v.a.c. X . En este tipo de variables, a diferencia de las discretas, la probabilidad de que la variable tome un valor aislado cualquiera es nula, y sólo hablaremos de probabilidades asociadas a intervalos.

Función de densidad

La *distribución de probabilidad* de X se suele caracterizar mediante una función $f(x)$, conocida como *función de densidad*. Formalmente, una función de densidad es una función no negativa, integrable en \mathbb{R} , que cumple

$$\int_{-\infty}^{\infty} f(x)dx = 1$$

A partir de esta función, se puede calcular la probabilidad de que el valor de la variable pertenezca a un intervalo $[a, b]$, midiendo el área encerrada por dicha función y el eje de abscisas entre los límites del intervalo, como se observa en la figura 8.1, es decir

$$P(a \leq X \leq b) = \int_a^b f(x)dx$$

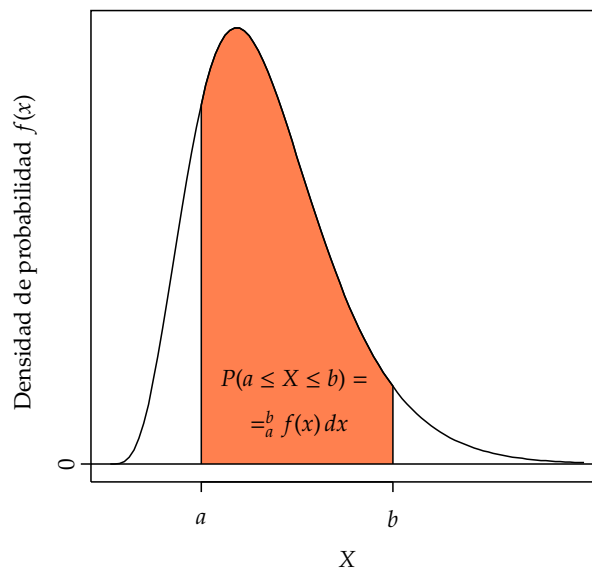


Figura 8.1 – En una v.a.c. la probabilidad asociada a un intervalo $[a, b]$, es el área que queda encerrada por la función de densidad y el eje de abscisas entre los límites del intervalo.

Función de distribución

Otra forma equivalente de caracterizar la distribución de probabilidad de X es mediante otra función $F(x)$, llamada *función de distribución*, que asigna a cada $x \in \mathbb{R}$ la probabilidad de que X tome un valor menor o igual que dicho número x . Así,

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(t) dt$$

A partir de la definición anterior es claro que la probabilidad de que la variable tome un valor en el intervalo $[a, b]$ puede calcularse a partir de la función de distribución de la siguiente forma:

$$P(a \leq X \leq b) = \int_a^b f(x) dx = \int_{-\infty}^b f(x) dx - \int_{-\infty}^a f(x) dx = F(b) - F(a)$$

Estadísticos poblacionales

Los parámetros descriptivos más importantes de una v.a.c. X son:

Media o Esperanza

$$E[X] = \mu = \int_{-\infty}^{\infty} x f(x) dx$$

Varianza

$$V[X] = \sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2$$

Desviación típica

$$D[X] = \sigma = +\sqrt{\sigma^2}$$

La media es una medida de tendencia central, mientras que la varianza y la desviación típica son medidas de dispersión.

Distribución Uniforme Continua

Una v.a.c. X se dice que sigue una *Distribución Uniforme Continua* de parámetros a y b , y se designa por $X \sim U(a, b)$, si su recorrido es el intervalo $[a, b]$ y su función de densidad es

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{si } a \leq x \leq b, \\ 0 & \text{en el resto} \end{cases}$$

Esta función es constante en el intervalo $[a, b]$ y nula fuera de él. Se cumple que

$$\mu = \frac{a+b}{2} \quad \sigma = +\frac{b-a}{\sqrt{12}}.$$

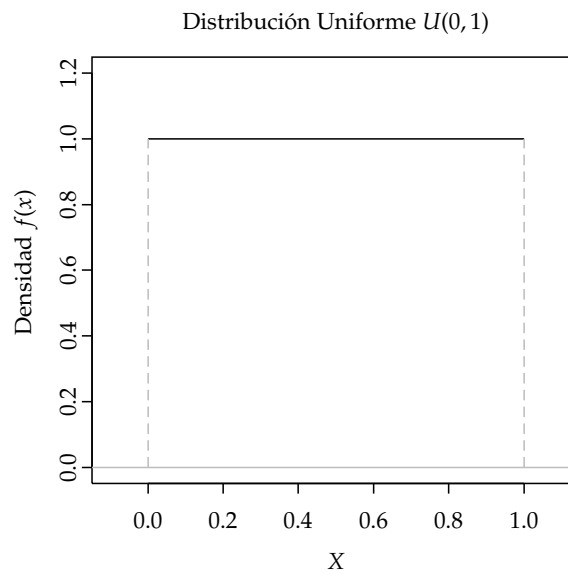


Figura 8.2 – Función de densidad de una variable aleatoria uniforme continua $U(0, 1)$.

Distribución Normal

Una v.a.c. X se dice que sigue una *Distribución Normal* o *Gaussiana* de media μ y desviación típica σ , y se designa por $X \sim N(\mu, \sigma)$, si su recorrido es todo \mathbb{R} y su función de densidad es

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Esta función tiene forma acampanada y es simétrica con respecto a la media μ .

La distribución Normal es la distribución continua más importante, ya que muchos de los fenómenos que aparecen en la naturaleza presentan esta distribución. Ello es debido a que, como establece el *Teorema Central del Límite*, cuando los resultados de un experimento están influidos por muchas causas independientes que actúan sumando sus efectos, se puede esperar que dichos resultados sigan una distribución normal.

La v.a.c normal de media 0 y desviación típica 1, $Z \sim N(0, 1)$, se conoce como *variable normal estándar* o *tipificada* y se utiliza muy a menudo. Su función de densidad aparece en la figura 8.3(a) y su función de distribución en la figura 8.3(b).

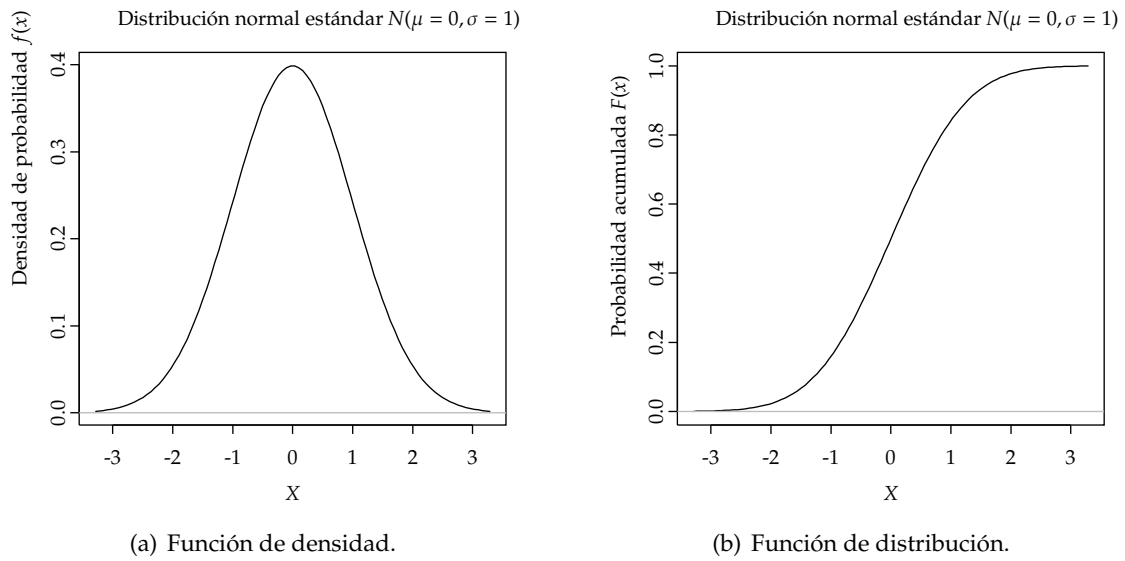


Figura 8.3 – Función de densidad y función de distribución de la variable aleatoria continua Z Normal de media 0 y desviación típica 1 $Z \sim N(0, 1)$

Distribución Chi-cuadrado

Si Z_1, \dots, Z_n son n v.a.c. normales estándar independientes, entonces la variable

$$X = Z_1^2 + \dots + Z_n^2$$

se dice que sigue una distribución *Chi-cuadrado* con n grados de libertad, y se nota $X \sim \chi^2(n)$.

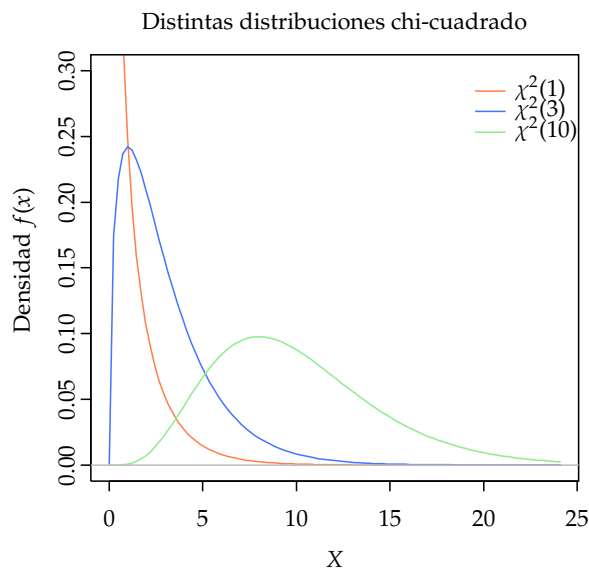


Figura 8.4 – Función de densidad de una variable aleatoria Chi cuadrado de 6 grados de libertad

Se cumple que

$$\begin{aligned}\mu &= n \\ \sigma &= +\sqrt{2n}\end{aligned}$$

La distribución Chi-cuadrado se utiliza en inferencia estadística para cálculos de intervalos de confianza y contrastes de hipótesis sobre la varianza de la población.

Distribución T de Student

Si Z es una v.a.c. normal estándar y X es una v.a.c chi-cuadrado con n grados de libertad, ambas variables independientes, entonces la variable

$$T = \frac{Z}{\sqrt{X/n}}$$

se dice que sigue una distribución T de Student con n grados de libertad, y se nota $T \sim T(n)$.

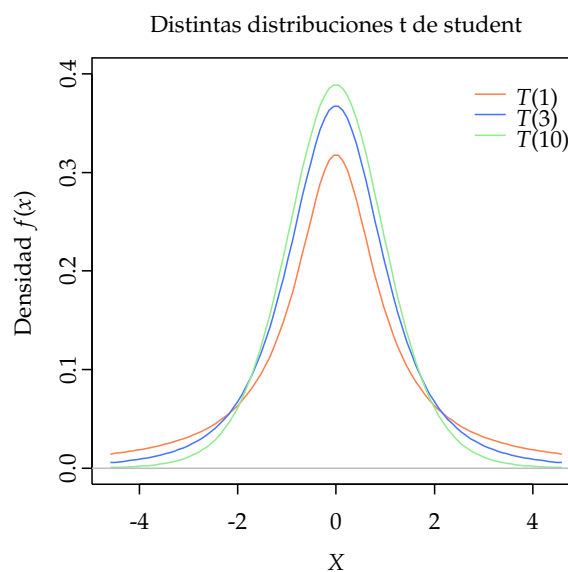


Figura 8.5 – Función de densidad de una variable aleatoria t de student de 10 grados de libertad

Esta variable es muy parecida a la normal estándar pero un poco menos apuntada, y se parece más a ésta a medida que aumentan los grados de libertad, de manera que para $n \geq 30$ ambas distribuciones se consideran prácticamente iguales. Se cumple que

$$\begin{aligned}\mu &= 0 \\ \sigma &= +\sqrt{n/(n-2)} \quad \text{con } n > 2\end{aligned}$$

La distribución T de Student se utiliza en inferencia estadística para cálculos de intervalos de confianza y contrastes de hipótesis sobre la media de la población.

Distribución F de Fisher-Snedecor

Si X e Y son dos v.a.c chi-cuadrado con m y n grados de libertad respectivamente, ambas variables independientes, entonces la variable

$$F = \frac{X/m}{Y/n}$$

se dice que sigue una distribución F de Fisher-Snedecor con m y n grados de libertad, y se denota $F \sim F(m, n)$.

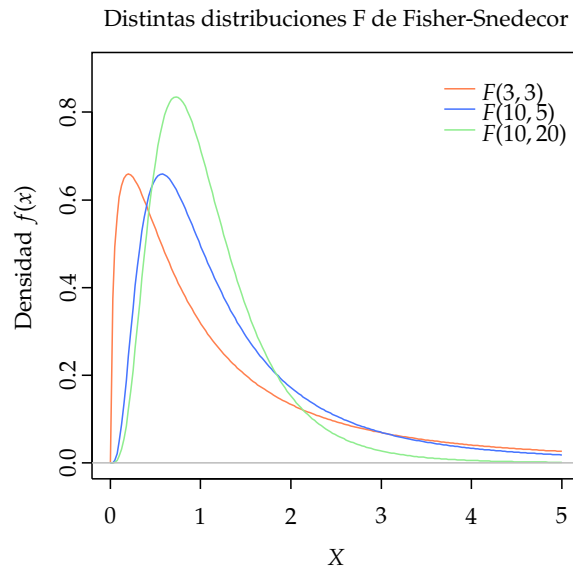


Figura 8.6 – Función de densidad de una variable aleatoria F de Fisher-Snedecor de 6 y 8 grados de libertad

$$\mu = \frac{n}{n-2}$$

$$\sigma = + \sqrt{\frac{2n^2(m+n-2)}{m(n-2)^2(n-4)}} \quad \text{con } n > 4$$

De la definición se deduce fácilmente que $F(m, n) = \frac{1}{F(n, m)}$, y si llamamos $F(m, n)_p$ al valor que cumple que $P(F(m, n) \leq F(m, n)_p) = p$, entonces se verifica

$$F(m, n)_p = \frac{1}{F(n, m)_{1-p}}$$

La distribución F de Fisher-Snedecor se utiliza en inferencia estadística para cálculos de intervalos de confianza y contrastes de hipótesis sobre el cociente de varianzas de dos poblaciones, y en análisis de la varianza.

2 Ejercicios resueltos

1. Supongase que un autobús pasa por una parada cada 15 minutos y que una persona puede llegar a la parada en cualquier instante, entonces la variable que mide el tiempo que la persona espera al autobús es una variable Uniforme continua $U(0, 15)$, ya que cualquier valor entre 0 y 15 minutos es equiprobable. Se pide:

a) Dibujar la gráfica de la función de densidad de la Uniforme $X \sim U(0, 15)$.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ Uniforme ▶ Gráfico de probabilidad.
- 2) En el cuadro de diálogo que aparece, introducir el valor 0 en el campo Mínimo, 15 en el campo Máximo y hacer click en el botón Aceptar.

b) Dibujar la gráfica de la función de distribución.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ Uniforme ▶ Gráfico de probabilidad.
- 2) En el cuadro de diálogo que aparece, introducir el valor 0 en el campo Mínimo, 15 en el campo Máximo, marcar la opción Función de distribución y hacer click en el botón Aceptar.

c) Calcular la probabilidad de esperar al autobús menos de 5 minutos.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ Uniforme ▶ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir el valor 5 en el campo Valores de la variable, 0 en el campo Mínimo, 15 en el campo Máximo y hacer click en el botón Aceptar.

d) Calcular la probabilidad de esperar al autobús más de 12 minutos.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ Uniforme ▶ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir el valor 12 en el campo Valores de la variable, 0 en el campo Mínimo, 15 en el campo Máximo, seleccionar la opción derecha en el campo cola de acumulación y hacer click en el botón Aceptar.

e) Calcular la probabilidad de esperar al autobús entre 5 y 10 minutos.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ Uniforme ▶ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir los valores 10, 5 en el campo Valores de la variable, 0 en el campo Mínimo, 15 en el campo Máximo y hacer click en el botón Aceptar.

La probabilidad del intervalo $P(5 \leq X \leq 10)$ es la resta de las probabilidades obtenidas $P(X \leq 10) - P(X \leq 5)$.

f) ¿Por debajo de qué tiempo esperará al autobús la mitad de las veces?



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ Uniforme ▶ Cuantiles.

- 2) En el cuadro de diálogo que aparece, introducir la probabilidad 0.5 en el campo Probabilidades acumuladas, 0 en el campo Mínimo, 15 en el campo Máximo y hacer click en el botón Aceptar.

g) ¿Por encima de qué tiempo esperará al autobús el 10 % de las veces?



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Uniforme ▸ Cuantiles.
- 2) En el cuadro de diálogo que aparece, introducir la probabilidad 0.1 en el campo Probabilidades acumuladas, 0 en el campo Mínimo, 15 en el campo Máximo, seleccionar la opción derecha en el campo cola de acumulación y hacer click en el botón Aceptar.

2. La variable aleatoria normal de media 0 y desviación típica 1, $Z \sim N(0, 1)$, se conoce como normal estándar y es la variable normal más importante. Se pide:

a) Dibujar la gráfica de la función de densidad.



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Normal ▸ Gráfico de probabilidad.
- 2) En el cuadro de diálogo que aparece, introducir el valor 0 en el campo Media, 1 en el campo Desviación típica y hacer click en el botón Aceptar.

b) ¿Cómo afectan los dos parámetros de la normal, su media y su desviación típica, a la forma de la campana de Gauss?



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Normal ▸ Gráfico de probabilidad.
- 2) En el cuadro de diálogo que aparece, seleccionar la opción Previsualizar.
- 3) Incrementar el valor de la media y ver cómo cambia la forma de la campana.
- 4) Después disminuir el valor de la desviación típica y ver cómo cambia la forma de la campana.

c) Dibujar la gráfica de la función de distribución.



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Normal ▸ Gráfico de probabilidad.
- 2) En el cuadro de diálogo que aparece, introducir el valor 0 en el campo Media, 1 en el campo Desviación típica, marcar la opción Función de distribución y hacer click en el botón Aceptar.

d) Calcular la probabilidad de que la normal estándar tome un valor menor que -1 .



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Normal ▸ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir el valor -1 en el campo Valores de la variable, 0 en el campo Media, 1 en el campo Desviación típica, y hacer click en el botón Aceptar.

e) Calcular la probabilidad de que la normal estándar tome un valor mayor que 1.



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Normal ▸ Probabilidades.

- 2) En el cuadro de diálogo que aparece, introducir el valor 1 en el campo Valores de la variable, 0 en el campo Media, 1 en el campo Desviación típica, seleccionar la opción derecha en el campo cola de acumulación y hacer click en el botón Aceptar.

f) Calcular la probabilidad de que la normal estándar tome un valor entre -1 (la media menos la desviación típica) y 1 (la media más la desviación típica).



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Normal ▸ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir los valores 1, -1 en el campo Valores de la variable, 0 en el campo Media, 1 en el campo Desviación típica y hacer click en el botón Aceptar.

La probabilidad del intervalo $P(-1 \leq Z \leq 1)$ es la resta de las probabilidades obtenidas $P(Z \leq 1) - P(Z \leq -1)$.

g) Calcular la probabilidad de que la normal estándar tome un valor entre -2 (la media menos dos veces la desviación típica) y 2 (la media más dos veces la desviación típica).



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Normal ▸ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir los valores 2, -2 en el campo Valores de la variable, 0 en el campo Media, 1 en el campo Desviación típica y hacer click en el botón Aceptar.

La probabilidad del intervalo $P(-2 \leq Z \leq 2)$ es la resta de las probabilidades obtenidas $P(Z \leq 2) - P(Z \leq -2)$.

h) Calcular la probabilidad de que la normal estándar tome un valor entre -3 (la media menos tres veces la desviación típica) y 3 (la media más tres veces la desviación típica).



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Normal ▸ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir los valores 3, -3 en el campo Valores de la variable, 0 en el campo Media, 1 en el campo Desviación típica y hacer click en el botón Aceptar.

La probabilidad del intervalo $P(-3 \leq Z \leq 3)$ es la resta de las probabilidades obtenidas $P(Z \leq 3) - P(Z \leq -3)$.

i) Calcular los cuantiles.



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Normal ▸ Cuantiles.
- 2) En el cuadro de diálogo que aparece, introducir las probabilidades 0.25, 0.5, 0.75 en el campo Probabilidades acumuladas, 0 en el campo Media, 1 en el campo Desviación típica y hacer click en el botón Aceptar.

j) Calcular el valor que deja acumulada por debajo una probabilidad 0,95.



- 1) Seleccionar el menú Teaching ▸ Distribuciones ▸ Continuas ▸ Normal ▸ Cuantiles.
- 2) En el cuadro de diálogo que aparece, introducir la probabilidad 0.95 en el campo Probabilidades acumuladas, 0 en el campo Media, 1 en el campo Desviación típica y hacer click en el botón Aceptar.

k) Calcular el valor que deja acumulada por encima una probabilidad 0,025.



- 1) Seleccionar el menú Teaching►Distribuciones►Continuas►Normal►Cuantiles.
- 2) En el cuadro de diálogo que aparece, introducir la probabilidad 0.025 en el campo Probabilidades acumuladas, 0 en el campo Media, 1 en el campo Desviación típica, seleccionar la opción derecha en el campo cola de acumulación y hacer click en el botón Aceptar.

3. El teorema central del límite establece que la variable resultante de sumar 30 o más variables independientes sigue una distribución normal de media la suma de las medias de cada una de las variables y de varianza la suma de sus varianzas. Esta es la explicación de que una gran parte de las variables continuas que aparecen en la naturaleza sean variables normales. Para observar de manera experimental el teorema central del límite se realiza un experimento que consiste en lanzar varios dados muchas veces y sumar los valores obtenidos. Se pide:

- a) Simular el lanzamiento de un dado 100000 veces y dibujar el diagrama de barras asociado. ¿Tiene forma de campana de Gauss?



Para generar los lanzamientos del dado:

- 1) Seleccionar el menú Teaching►Simulaciones►Lanzamiento de dados.
- 2) En el cuadro de diálogo que aparece, introducir 1 en el campo Número de dados, introducir 100000 en el campo Número de lanzamientos, seleccionar la opción Incluir suma, introducir un nombre para el conjunto de datos y hacer click en el botón Aceptar.

Para dibujar el diagrama de barras:

- 1) Seleccionar el menú Teaching►Gráficos►Diagrama de barras.
- 2) En el cuadro de diálogo que aparece seleccionar la variable sum.
- 3) En la solapa Opciones de las barras, seleccionar la opción Frecuencias relativas y hacer click en el botón Aceptar.

- b) Repetir el apartado anterior con 2 y 30 dados. ¿Se cumple el teorema central del límite?

4. La suma de n variables normales estándar independientes elevadas al cuadrado es una variable con distribución Chi-cuadrado con n grados de libertad $\chi^2(n)$. Sea X una variable Chi-cuadrado con 6 grados de libertad $\chi^2(6)$. Se pide:

- a) Dibujar la gráfica de la función de densidad.



- 1) Seleccionar el menú Teaching►Distribuciones►Continuas►Chi-cuadrado►Gráfico de probabilidad.
- 2) En el cuadro de diálogo que aparece, introducir el valor 6 en el Grados de libertad y hacer click en el botón Aceptar.

- b) Calcular la probabilidad de que la variable tome un valor menor que 6.



- 1) Seleccionar el menú Teaching►Distribuciones►Continuas►Chi-cuadrado►Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir el valor 6 en el campo Valores de la variable, 6 en el campo Grados de libertad y hacer click en el botón Aceptar.

- c) Calcular el valor que deja acumulada por debajo una probabilidad 0,05.



- 1) Seleccionar el menú Teaching►Distribuciones►Continuas►Chi-cuadrado►Cuantiles.

- 2) En el cuadro de diálogo que aparece, introducir la probabilidad 0.05 en el campo Probabilidades acumuladas, 6 en el campo Grados de libertad y hacer click en el botón Aceptar.

d) Calcular el valor que deja acumulada por arriba una probabilidad 0,1.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ Chi-cuadrado ▶ Cuantiles.
- 2) En el cuadro de diálogo que aparece, introducir la probabilidad 0.1 en el campo Probabilidades, 6 en el campo Grados de libertad, seleccionar la opción derecha en el campo cola de acumulación y hacer click en el botón Aceptar.

5. La variable que se obtiene al dividir una normal estándar entre la raíz de una variable Chi-cuadrado de n grados de libertad dividida por n , sigue una distribución t de student de n grados de libertad $T(n)$. Sea X una variable t de student de 8 grados de libertad $T(8)$. Se pide:

a) Dibujar la gráfica de la función de probabilidad y compararla con la de la normal estándar.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ T de student ▶ Gráfico de probabilidad.
- 2) En el cuadro de diálogo que aparece, introducir el valor 8 en el campo Grados de libertad y hacer click en el botón Aceptar.

b) Calcular el percentil octavo.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ T de student ▶ Cuantiles.
- 2) En el cuadro de diálogo que aparece, introducir la probabilidad 0.08 en el campo Probabilidades acumuladas, 8 en el campo Grados de libertad y hacer click en el botón Aceptar.

c) Calcular el valor por encima del cual está el 5 % de la población.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ T de student ▶ Cuantiles.
- 2) En el cuadro de diálogo que aparece, introducir la probabilidad 0.05 en el campo Probabilidades, 8 en el campo Grados de libertad, seleccionar la opción derecha en el campo cola de acumulación y hacer click en el botón Aceptar.

6. La variable resultante de dividir una variable Chi-cuadrado de n grados de libertad dividida por n , entre una variable Chi-cuadrado de m grados de libertad dividida por m , sigue un modelo de distribución F de Fisher de n y m grados de libertad $F(n, m)$. Sea X una variable F de Fisher de 10 y 20 grados de libertad $F(10, 20)$. Se pide:

a) Dibujar la gráfica de la función de densidad



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ F de Fisher ▶ Gráfico de probabilidad.
- 2) En el cuadro de diálogo que aparece, introducir 10 el campo Grados de libertad del numerador, introducir 20 en el campo Grados de libertad del denominador y hacer click en el botón Aceptar.

b) Calcular la probabilidad acumulada por encima de 1.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ F de Fisher ▶ Probabilidades.
- 2) En el cuadro de diálogo que aparece, introducir el valor 1 en el campo Valores de la variable, 10 en el campo Grados de libertad del numerador, 20 en el campo Grados de libertad del denominador, seleccionar la opción derecha en el campo cola de acumulación y hacer click en el botón Aceptar.

c) Calcular el rango intercuartílico.



- 1) Seleccionar el menú Teaching ▶ Distribuciones ▶ Continuas ▶ F de Fisher ▶ Cuantiles.
- 2) En el cuadro de diálogo que aparece, introducir las probabilidades 0.75, 0.25 en el campo Probabilidades, 10 en el campo Grados de libertad del numerador, 20 en el campo Grados de libertad del denominador y hacer click en el botón Aceptar.

El rango intercuartílico es la resta de los valores obtenidos correspondientes al tercer y primer cuantiles.

3 Ejercicios propuestos

1. Entre los diabéticos, el nivel de glucosa en la sangre en ayunas X , puede suponerse de distribución aproximadamente normal, con media 106mg/100ml y desviación típica 8mg/100ml.
 - a) Hallar $P(X \leq 120\text{mg}/100\text{ml})$
 - b) ¿Qué porcentaje de diabéticos tendrá niveles entre 90 y 120mg/100ml?
 - c) Encontrar un valor que tenga la propiedad de que el 25 % de los diabéticos tenga un nivel de glucosa por debajo de dicho valor.
2. Se sabe que el nivel de colesterol en varones de más de 30 años de una determinada población sigue una distribución normal, de media 220mg/dl y desviación típica 30mg/dl. Si la población tiene 20000 varones mayores de 30 años,
 - a) ¿Cuántos se espera que tengan su nivel de colesterol entre 210mg/dl y 240mg/dl?
 - b) ¿Cuántos se espera que tengan su nivel de colesterol por encima de 250mg/dl?
 - c) ¿Cuál será el nivel de colesterol por encima del cual se espera que esté el 20 % de la población?
3. Calcular la probabilidad de obtener entre 40 y 60 caras, inclusive, al lanzar 100 veces una moneda. Utilizar la aproximación de la distribución binomial mediante una normal.

Intervalos de Confianza para Medias y Proporciones

1 Fundamentos teóricos

1.1 Inferencia Estadística y Estimación de Parámetros

El objetivo de un estudio estadístico es doble: describir la muestra elegida de una población en la que se quiere estudiar alguna característica, y realizar inferencias, es decir, sacar conclusiones y hacer predicciones sobre la población de la que se ha extraído dicha muestra.

La metodología que conduce a obtener conclusiones sobre la población, basadas en la información contenida en la muestra, constituye la *Inferencia Estadística*.

Puesto que la muestra contiene menos información que la población, las predicciones serán aproximadas. Por eso, uno de los objetivos de la inferencia estadística es determinar la probabilidad de que una conclusión obtenida a partir del análisis de una muestra sea cierta, y para ello se apoya en la teoría de la probabilidad.

Cuando se desea conocer el valor de alguno de los parámetros de la población, el procedimiento a utilizar es la *Estimación de Parámetros*, que a su vez se divide en *Estimación Puntual*, cuando se da un único valor como estimación del parámetro poblacional considerado, y *Estimación por Intervalos*, cuando interesa conocer no sólo un valor aproximado del parámetro sino también la precisión de la estimación. En este último caso el resultado es un intervalo, dentro del cual estará, con una cierta confianza, el verdadero valor del parámetro poblacional. A este intervalo se le denomina *intervalo de confianza*. A diferencia de la estimación puntual, en la que se utiliza un único estimador, en la estimación por intervalo se emplean dos estimadores, uno para cada extremo del intervalo.

1.2 Intervalos de Confianza

Dados dos estadísticos muestrales L_1 y L_2 , se dice que el intervalo $I = (L_1, L_2)$ es un *Intervalo de Confianza* para un parámetro poblacional θ , con *nivel de confianza* $1 - \alpha$ (o *nivel de significación* α), si la probabilidad de que los estadísticos que determinan los límites del intervalo tomen valores tales que θ esté comprendido entre ellos, es igual a $1 - \alpha$, es decir,

$$P(L_1 < \theta < L_2) = 1 - \alpha$$

Los extremos del intervalo son variables aleatorias cuyos valores dependen de la muestra considerada. Es decir, los extremos inferior y superior del intervalo serían $L_1(X_1, \dots, X_n)$ y $L_2(X_1, \dots, X_n)$ respectivamente, aunque habitualmente escribiremos L_1 y L_2 para simplificar la notación. Designaremos mediante l_1 y l_2 los valores que toman dichas variables para una muestra determinada (x_1, \dots, x_n) .

Cuando en la definición se dice que la probabilidad de que el parámetro θ esté en el intervalo (L_1, L_2) es $1 - \alpha$, quiere decir que en el $100(1 - \alpha)\%$ de las posibles muestras, el valor de θ estaría en los correspondientes intervalos (l_1, l_2) .

Una vez que se tiene una muestra, y a partir de ella se determina el intervalo correspondiente (l_1, l_2) , no tendría sentido hablar de la probabilidad de que el parámetro θ esté en el intervalo (l_1, l_2) ,

pues al ser l_1 y l_2 números, el parámetro θ , que también es un número, aunque desconocido, estará o no estará en dicho intervalo, y por ello hablamos de confianza en lugar de probabilidad.

Así, cuando hablemos de un intervalo de confianza para el parámetro θ con nivel de confianza $1 - \alpha$, entenderemos que antes de tomar una muestra, hay una probabilidad $1 - \alpha$ de que el intervalo que se construya a partir de ella, contenga el valor del parámetro θ . O, dicho de otro modo, si tomásemos 100 muestras del mismo tamaño y calculásemos sus respectivos intervalos, el $1 - \alpha$ % de estos contendrían el verdadero valor del parámetro a estimar (ver figura 9.1).

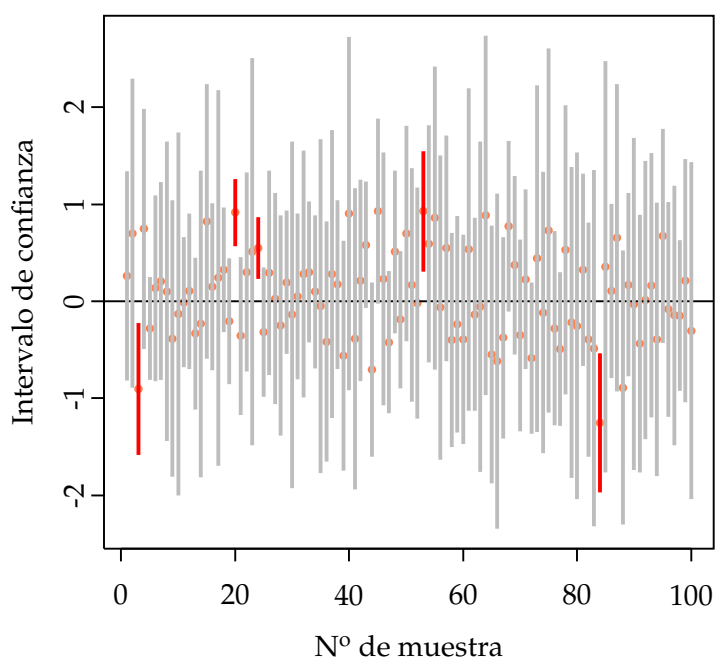


Figura 9.1 – Intervalos de confianza del 95 % para la media de 100 muestras tomadas de una población normal $N(0, 1)$. Como se puede apreciar, de los 100 intervalos, sólo 5 no contienen el valor de la media real $\mu = 0$.

Cuando se realiza la estimación de un parámetro mediante un intervalo de confianza, el nivel de confianza se suele fijar a niveles altos (los más habituales son 0,90, 0,95 ó 0,99), para tener una alta confianza de que el parámetro está dentro del intervalo. Por otro lado, también interesa que la amplitud del intervalo sea pequeña para delimitar con precisión el valor del parámetro poblacional (esta amplitud del intervalo se conoce como *imprecisión* de la estimación). Pero a partir de una muestra, cuanto mayor sea el nivel de confianza deseado, mayor amplitud tendrá el intervalo y mayor imprecisión la estimación, y si se impone que la estimación sea más precisa (menor imprecisión), el nivel de confianza correspondiente será más pequeño. Por consiguiente, hay que llegar a una solución de compromiso entre el nivel de confianza y la precisión de la estimación. No obstante, si con la muestra disponible no es posible obtener un intervalo de amplitud suficientemente pequeña (imprecisión pequeña) con un nivel de confianza aceptable, hay que emplear una muestra de mayor tamaño. Al aumentar el tamaño muestral se consiguen intervalos de menor amplitud sin disminuir el nivel de confianza, o niveles de confianza más altos manteniendo la amplitud.

Intervalos de confianza para la media

Apoyándose en conclusiones extraídas del Teorema Central del Límite se obtiene que, siempre que las muestras sean grandes (como criterio habitual se toma que el tamaño muestral, n , sea mayor o igual que 30), e independientemente de la distribución original de la variable de partida X , de media

9. Intervalos de Confianza para Medias y Proporciones

μ y desviación típica σ , la variable

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

sigue una distribución Normal tipificada, $N(0, 1)$.

Si la desviación típica σ de la variable de partida es desconocida, se utiliza como estimación la cuasidesviación típica muestral:

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

y con ello, la nueva variable

$$\frac{\bar{X} - \mu}{\hat{S} / \sqrt{n}}$$

sigue una distribución t de Student con $n - 1$ grados de libertad, $T(n - 1)$.

Para muestras pequeñas ($n < 30$) también pueden aplicarse los resultados anteriores, siempre y cuando la variable aleatoria de partida X , siga una distribución Normal.

A partir de lo anterior y teniendo en cuenta los tres factores de clasificación expuestos: si la población de partida en la que obtenemos la muestra sigue o no una distribución Normal, si la varianza de dicha población es conocida o desconocida, y si la muestra es grande ($n \geq 30$) o no, pueden deducirse las siguientes expresiones correspondientes a los diferentes intervalos de confianza.

Intervalo de confianza para la media de una población normal con varianza conocida en muestras de cualquier tamaño

$$\left(\bar{x} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right)$$

En la figura 9.2 aparece un esquema explicativo de la construcción de este intervalo.

Intervalo de confianza para la media de una población normal con varianza desconocida en muestras de cualquier tamaño

$$\left(\bar{x} - t_{\alpha/2}^{n-1} \cdot \frac{\hat{s}}{\sqrt{n}}, \bar{x} + t_{\alpha/2}^{n-1} \cdot \frac{\hat{s}}{\sqrt{n}} \right)$$

Si las muestras son grandes ($n \geq 30$) el anterior intervalo puede aproximarse mediante:

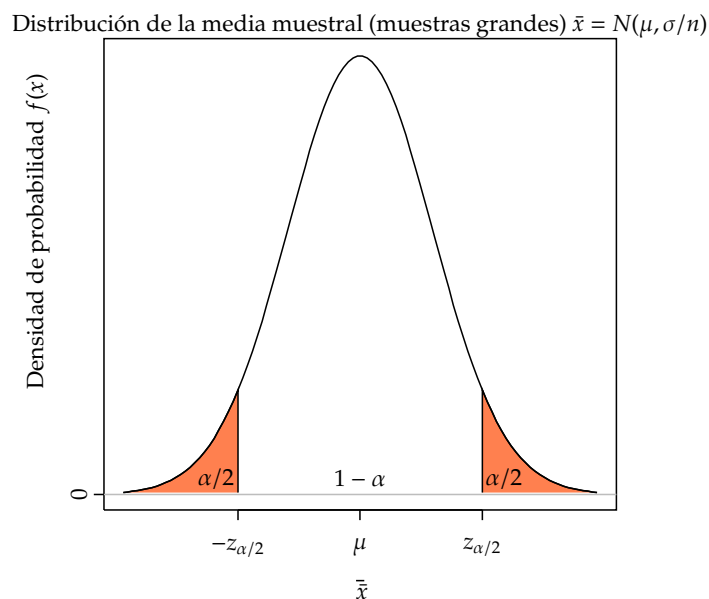
$$\left(\bar{x} - z_{\alpha/2} \cdot \frac{\hat{s}}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \cdot \frac{\hat{s}}{\sqrt{n}} \right)$$

Intervalo de confianza para la media de una población no normal, varianza conocida y muestras grandes ($n \geq 30$)

$$\left(\bar{x} - z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \right)$$

Intervalo de confianza para la media de una población no normal, varianza desconocida y muestras grandes ($n \geq 30$)

$$\left(\bar{x} - t_{\alpha/2}^{n-1} \cdot \frac{\hat{s}}{\sqrt{n}}, \bar{x} + t_{\alpha/2}^{n-1} \cdot \frac{\hat{s}}{\sqrt{n}} \right)$$



$$P\left(\mu - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \bar{x} \leq \mu + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

Figura 9.2 – Cálculo del intervalo de confianza para la media de una población normal con varianza conocida, a partir de la distribución de la media muestral $\bar{x} \sim N(\mu, \sigma/\sqrt{n})$ para muestras grandes ($n \geq 30$).

Al tratarse de muestras grandes, el anterior intervalo puede aproximarse por:

$$\left(\bar{x} - z_{\alpha/2} \cdot \frac{\hat{s}}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \cdot \frac{\hat{s}}{\sqrt{n}}\right)$$

Si la población de partida no es normal, y las muestras son pequeñas, no puede aplicarse el Teorema Central del Límite y no se obtienen intervalos de confianza para la media.

Para cualquiera de los anteriores intervalos:

- n es el tamaño de la muestra.
- \bar{x} es la media muestral.
- σ es la desviación típica de la población.
- \hat{s} es la cuasidesviación típica muestral: $\hat{s}^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1}$.
- $z_{\alpha/2}$ es el valor que deja a su derecha una probabilidad $\alpha/2$ en una distribución Normal tipificada.
- $t_{\alpha/2}^{n-1}$ es el valor que deja a su derecha una probabilidad $\alpha/2$ en una distribución t de Student con $n - 1$ grados de libertad.

Intervalos de confianza para la proporción poblacional p

Para muestras grandes ($n \geq 30$) y valores de p (probabilidad de “éxito”) cercanos a 0,5, la distribución Binomial puede aproximarse mediante una Normal de media np y desviación típica $\sqrt{np(1-p)}$. En la práctica, para que sea válida dicha aproximación, se toma el criterio de que tanto np como $n(1-p)$ deben ser mayores que 5. Esto hace que también podamos construir intervalos de confianza para proporciones tomando éstas como medias de variables dicotómicas en las que la presencia o ausencia

9. Intervalos de Confianza para Medias y Proporciones

de la característica objeto de estudio (“éxito” ó “fracaso”) se expresan mediante un 1 ó un 0 respectivamente.

De este modo, en muestras grandes y con distribuciones binomiales no excesivamente asimétricas (tanto np como $n(1 - p)$ deben ser mayores que 5), si denominamos \widehat{p} a la proporción de individuos que presentan el atributo estudiado en la muestra concreta, entonces el intervalo de confianza para la proporción con un nivel de significación α viene dado por:

$$\left(\widehat{p} - z_{\alpha/2} \cdot \sqrt{\frac{\widehat{p} \cdot (1 - \widehat{p})}{n}}, \widehat{p} + z_{\alpha/2} \cdot \sqrt{\frac{\widehat{p} \cdot (1 - \widehat{p})}{n}} \right)$$

donde:

- n es el tamaño muestral.
- \widehat{p} a la proporción de individuos que presentan el atributo estudiado en la muestra concreta.
- $z_{\alpha/2}$ es el valor que deja a su derecha una probabilidad $\alpha/2$ en una distribución Normal tipificada.

En muestras pequeñas o procedentes de una Binomial fuertemente asimétrica ($np \leq 5$ ó $n(1-p) \leq 5$) no puede aplicarse el Teorema Central del Límite y la construcción de intervalos de confianza debe realizarse a partir de la distribución Binomial.

2 Ejercicios resueltos

1. Se analiza la concentración de principio activo en una muestra de 10 envases tomados de un lote de un fármaco, obteniendo los siguientes resultados en mg/mm^3 :

17,6 – 19,2 – 21,3 – 15,1 – 17,6 – 18,9 – 16,2 – 18,3 – 19,0 – 16,4

Se pide:

- Crear un conjunto de datos con la variable concentracion.
- Calcular el intervalo de confianza para la media de la concentración del lote con nivel de confianza del 95 % (nivel de significación $\alpha = 0,05$).



- Seleccionar el menú Teaching ▶ Test paramétricos ▶ Medias ▶ Test t para una muestra.
- En el cuadro de diálogo que aparece seleccionar la variable concentracion en el campo Variable y hacer clic sobre el botón Enviar.

- Calcular los intervalos de confianza para la media con niveles del 90 % y del 99 % (niveles de significación $\alpha = 0,1$ y $\alpha = 0,01$).



Repetir los mismos pasos del apartado anterior, cambiando el nivel de confianza para cada intervalo en la solapa Opciones de contraste

- Si definimos la precisión del intervalo como la inversa de su amplitud, ¿cómo afecta a la precisión del intervalo de confianza el tomar niveles de significación cada vez más altos? ¿Cuál puede ser la explicación?
- ¿Qué tamaño muestral sería necesario para obtener una estimación del contenido medio de principio activo con un margen de error de $\pm 0,5 \text{ mg}/\text{mm}^3$ y una confianza del 95 %?



- Seleccionar el menú Teaching ▶ Estadística descriptiva ▶ Estadísticos.
- En el cuadro de diálogo que aparece seleccionar la variable concentracion en el campo Variable.
- En la solapa Estadísticos básicos marcar el estadístico Cuasidesviación típica y hacer clic en el botón Enviar.
- Seleccionar el menú Teaching ▶ Test paramétricos ▶ Medias ▶ Cálculo del tamaño muestral para la media.
- En el cuadro de diálogo que aparece introducir la cuasidesviación típica muestral en el campo Desviación típica, el nivel de confianza deseado, en este caso 0,95, en el campo Nivel de confianza, el margen de error deseado, en este caso 0,5, en el campo Error, y hacer clic en el botón Enviar.

- Si, para que sea efectivo, el fármaco debe tener una concentración mínima de $16 \text{ mg}/\text{mm}^3$ de principio activo, ¿se puede aceptar el lote como bueno? Justificar la respuesta.

2. Una central de productos lácteos recibe diariamente la leche de dos granjas X e Y. Para analizar la calidad de la leche, durante una temporada, se controla el contenido de materia grasa de la leche

9. Intervalos de Confianza para Medias y Proporciones

que proviene de ambas granjas, con los siguientes resultados:

X		Y	
0,34	0,34	0,28	0,29
0,32	0,35	0,30	0,32
0,33	0,33	0,32	0,31
0,32	0,32	0,29	0,29
0,33	0,30	0,31	0,32
0,31	0,32	0,29	0,31
		0,33	0,32
		0,32	0,33

- Crear un conjunto de datos con las variables grasa y granja.
- Calcular el intervalo de confianza con un 95 % de confianza para el contenido medio de materia grasa de la leche sin tener en cuenta si la misma procede de una u otra granja.



- Seleccionar el menú Teaching ▶ Test paramétricos ▶ Medias ▶ Test t para una muestra.
- En el cuadro de diálogo que aparece seleccionar la variable grasa en el campo Variable y hacer clic sobre el botón Enviar.

- Calcular los intervalos de confianza con un 95 % de confianza para el contenido medio de materia grasa de la leche dividiendo los datos según la granja de procedencia de la leche.



- Seleccionar el menú Teaching ▶ Test paramétricos ▶ Medias ▶ Test t para una muestra.
- En el cuadro de diálogo que aparece seleccionar la variable grasa en el campo Variable.
- Seleccionar la casilla de Filtro e introducir la condición granja== 'X' hacer clic sobre el botón Enviar.
- Repetir los mismos pasos para el intervalo de confianza de la granja Y, introduciendo la condición granja== 'Y' en el campo Condición de selección.

- A la vista de los intervalos obtenidos en el punto anterior, ¿se puede concluir que existen diferencias significativas en el contenido medio de grasa según la procedencia de la leche? Justificar la respuesta.

- En una encuesta realizada en una facultad, sobre si el alumnado utiliza habitualmente (al menos una vez a la semana) la biblioteca de la misma, se han obtenido los siguientes resultados:

Alumno	1	2	3	4	5	6	7	8	9	10	11	12
Respuesta	No	Si	No	No	No	Si	No	Si	Si	Si	Si	No

Alumno	13	14	15	16	17	18	19	20	21	22	23
Respuesta	Si	No	Si	No	No	No	Si	Si	Si	No	No

Alumno	24	25	26	27	28	29	30	31	32	33	34
Respuesta	Si	No	No	Si	Si	No	No	Si	No	Si	No

- Crear un conjunto de datos con la variable respuesta.
- Calcular el intervalo de confianza con $\alpha = 0,01$ para la proporción del alumnado que utiliza habitualmente la biblioteca.



- Seleccionar el menú Teaching ▶ Test paramétricos ▶ Proporciones ▶ Test para una proporción.

- 2) En el cuadro de diálogo que aparece seleccionar la variable respuesta en el campo Variable e introducir sí en el campo Categoría.
- 3) En la solapa Opciones de contraste introducir 0,99 en el campo Nivel de confianza y hacer clic en el botón Enviar.

- c) ¿Qué interpretación tiene dicho intervalo? ¿Cómo es su precisión?
- d) ¿Qué tamaño muestral sería necesario para obtener una estimación del porcentaje de alumnos que utilizan regularmente la biblioteca con un margen de error de un 1 % y una confianza del 95 %?



- 1) Seleccionar el menú Teaching ▶ Test paramétricos ▶ Proporciones ▶ Cálculo del tamaño muestral para una proporción.
- 2) En el cuadro de diálogo que aparece introducir la proporción muestral en el campo p , el nivel de confianza deseado, en este caso 0,95, en el campo Nivel de confianza, el margen de error deseado, en este caso 0,01, en el campo Error, y hacer clic en el botón Enviar.

4. El Ministerio de Sanidad está interesado en la elaboración de un intervalo de confianza para la proporción de personas mayores de 65 años con problemas respiratorios que han sido vacunadas en una determinada ciudad. Para ello, después de preguntar a 200 pacientes mayores de 65 años con problemas respiratorios en los hospitales de dicha ciudad, 154 responden afirmativamente.

- a) Calcular el intervalo de confianza al 95 % para la proporción de pacientes vacunados.



- 1) Seleccionar el menú Teaching ▶ Test paramétricos ▶ Proporciones ▶ Test para una proporción.
- 2) En el cuadro de diálogo que aparece marcar la opción Introducción manual de frecuencias, introducir 154 en el campo Frecuencia muestral, introducir 200 en el campo Tamaño muestral y hacer clic en el botón Enviar.

- b) Si entre los objetivos del Ministerio se encontraba alcanzar una proporción del al menos un 70 % de vacunados en dicho colectivo, ¿se puede concluir que se han cumplido los objetivos? Justificar la respuesta.

3 Ejercicios propuestos

1. Para determinar el nivel medio de colesterol (en mg/dl) en la sangre de una población, se realizaron análisis sobre una muestra de 8 personas, obteniéndose los siguientes resultados:

196 212 188 206 203 210 201 198

Hallar los intervalos de confianza para la media del nivel de colesterol con niveles de significación 0,1, 0,05 y 0,01. ¿Se puede afirmar que el nivel de colesterol medio de la población está por debajo de 210 mg/dl?

2. Para tratar un determinado síndrome neurológico se utilizan dos técnicas A y B . En un estudio se tomó una muestra de 60 pacientes con dicho síndrome y se le aplicó la técnica A a 25 de ellos y la técnica B a los 35 restantes. De los pacientes tratados con la técnica A , 18 se curaron, mientras que de los tratados con la técnica B , se curaron 21. Calcular un intervalo de confianza del 95 % para la proporción de curaciones con cada técnica. ¿Qué intervalo es más preciso?
3. A las siguientes elecciones locales en una ciudad se presentan tres partidos: A , B y C . Con el objetivo de hacer una estimación sobre la proporción de voto que cada uno de ellos obtendrá, se realiza una encuesta en la que responden 300 personas, de las cuales 60 piensan votar a A , 80 a

9. Intervalos de Confianza para Medias y Proporciones

- B, 90 a C, 15 en blanco y 55 abstenciones. Calcular un intervalo de confianza para la proporción de votos, sobre el total del censo, de cada uno de los partidos que se presentan.
4. El fichero `nations.txt` contiene información sobre el desarrollo de distintos países (tasa de uso de anticonceptivos (`contraception`), producto interior bruto per cápita (`GDP`), tasa de mortalidad infantil (`infant.mortality`) y tasa de fertilidad (`TFR`)). Se pide:
- Importar el fichero `nations.txt` en un conjunto de datos.
 - Calcular el intervalo de confianza de la tasa de uso de anticonceptivos y de la tasa de fertilidad para los países con un producto interior bruto per cápita superior a 10000 US\$ e inferiores a dicha cantidad. Interpretar los intervalos.

Intervalos de Confianza para la Comparación de 2 Poblaciones

1 Fundamentos teóricos

1.1 Inferencia Estadística y Estimación de Parámetros

El objetivo de un estudio estadístico es doble: describir la muestra elegida de una población en la que se quiere estudiar alguna característica, y realizar inferencias, es decir, sacar conclusiones y hacer predicciones sobre la población de la que se ha extraído dicha muestra.

La metodología que conduce a obtener conclusiones sobre la población, basadas en la información contenida en la muestra, constituye la *Inferencia Estadística*.

Puesto que la muestra contiene menos información que la población, las predicciones serán aproximadas. Por eso, uno de los objetivos de la inferencia estadística es determinar la probabilidad de que una conclusión obtenida a partir del análisis de una muestra sea cierta, y para ello se apoya en la teoría de la probabilidad.

Cuando se desea conocer el valor de alguno de los parámetros de la población, el procedimiento a utilizar es la *Estimación de Parámetros*, que a su vez se divide en *Estimación Puntual*, cuando se da un único valor como estimación del parámetro poblacional considerado, y *Estimación por Intervalos*, cuando interesa conocer no sólo un valor aproximado del parámetro sino también la precisión de la estimación. En este último caso el resultado es un intervalo, dentro del cual estará, con una cierta confianza, el verdadero valor del parámetro poblacional. A este intervalo se le denomina *intervalo de confianza*. A diferencia de la estimación puntual, en la que se utiliza un único estimador, en la estimación por intervalo emplearemos dos estimadores, uno para cada extremo del intervalo.

1.2 Intervalos de Confianza

Dados dos estadísticos muestrales L_1 y L_2 , se dice que el intervalo $I = (L_1, L_2)$ es un *Intervalo de Confianza* para un parámetro poblacional θ , con *nivel de confianza* $1 - \alpha$ (o *nivel de significación* α), si la probabilidad de que los estadísticos que determinan los límites del intervalo tomen valores tales que θ esté comprendido entre ellos, es igual a $1 - \alpha$, es decir,

$$P(L_1 < \theta < L_2) = 1 - \alpha$$

Los extremos del intervalo son variables aleatorias cuyos valores dependen de la muestra considerada. Es decir, los extremos inferior y superior del intervalo serían $L_1(X_1, \dots, X_n)$ y $L_2(X_1, \dots, X_n)$ respectivamente, aunque habitualmente escribiremos L_1 y L_2 para simplificar la notación. Designaremos mediante l_1 y l_2 los valores que toman dichas variables para una muestra determinada (x_1, \dots, x_n) .

Cuando en la definición se dice que la probabilidad de que el parámetro θ esté en el intervalo (L_1, L_2) es $1 - \alpha$, quiere decir que en el $100(1 - \alpha)\%$ de las posibles muestras, el valor de θ estaría en los correspondientes intervalos (l_1, l_2) .

Una vez que se tiene una muestra, y a partir de ella se determina el intervalo correspondiente (l_1, l_2) , no tendría sentido hablar de la probabilidad de que el parámetro θ esté en el intervalo (l_1, l_2) ,

pues al ser l_1 y l_2 números, el parámetro θ , que también es un número, aunque desconocido, estará o no estará en dicho intervalo, y por ello hablamos de confianza en lugar de probabilidad.

Así, cuando hablemos de un intervalo de confianza para el parámetro θ con nivel de confianza $1 - \alpha$, entenderemos que antes de tomar una muestra, hay una probabilidad $1 - \alpha$ de que el intervalo que se construya a partir de ella, contenga el valor del parámetro θ .

Cuando se realiza la estimación de un parámetro mediante un intervalo de confianza, el nivel de confianza se suele fijar a niveles altos (los más habituales son 0,90, 0,95 ó 0,99), para tener una alta confianza de que el parámetro está dentro del intervalo. Por otro lado, también interesa que la amplitud del intervalo sea pequeña para delimitar con precisión el valor del parámetro poblacional (esta amplitud del intervalo se conoce como *imprecisión* de la estimación). Pero a partir de una muestra, cuanto mayor sea el nivel de confianza deseado, mayor amplitud tendrá el intervalo y mayor imprecisión la estimación, y si se impone que la estimación sea más precisa (menor imprecisión), el nivel de confianza correspondiente será más pequeño. Por consiguiente, hay que llegar a una solución de compromiso entre el nivel de confianza y la precisión de la estimación. No obstante, si con la muestra disponible no es posible obtener un intervalo de amplitud suficientemente pequeña (imprecisión pequeña) con un nivel de confianza aceptable, hay que emplear una muestra de mayor tamaño. Al aumentar el tamaño muestral se consiguen intervalos de menor amplitud sin disminuir el nivel de confianza, o niveles de confianza más altos manteniendo la amplitud.

Intervalos de confianza para la diferencia de medias

De igual manera a como ocurría con los intervalos de confianza para la media de una variable, apoyándose en conclusiones extraídas del Teorema Central del Límite se puede demostrar que, en muestras grandes ($n_1 \geq 30$ y $n_2 \geq 30$), procedentes de poblaciones de dos variables X_1 y X_2 , con distribuciones no necesariamente Normales, de medias μ_1 y μ_2 y desviaciones típicas σ_1 y σ_2 respectivamente, la variable

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

sigue una distribución Normal tipificada, $N(0, 1)$.

De igual manera, si las varianzas de las variables son desconocidas, utilizando como estimadores muestrales sus correspondientes cuasivarianzas \hat{S}_1^2 y \hat{S}_2^2 , donde

$$\hat{S}_1^2 = \frac{\sum (x_{1,i} - \bar{x}_1)^2}{n_1 - 1} \quad \text{y} \quad \hat{S}_2^2 = \frac{\sum (x_{2,i} - \bar{x}_2)^2}{n_2 - 1}$$

entonces la variable

$$T = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\hat{S}_1^2}{n_1} + \frac{\hat{S}_2^2}{n_2}}}$$

sigue una distribución t de Student, en la que el número de grados de libertad dependerá de si las varianzas, aún siendo desconocidas, pueden considerarse iguales o no.

Para muestras pequeñas ($n_1 < 30$ ó $n_2 < 30$), las distribuciones anteriores son también aplicables siempre que las variables de partida sigan distribuciones Normales.

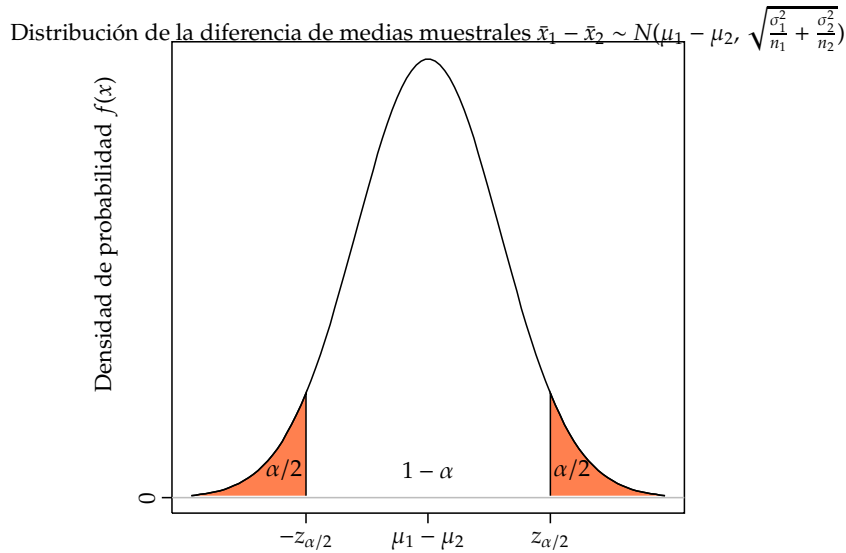
A partir de todo ello y teniendo en cuenta los tres factores de clasificación comentados: si las poblaciones de partida en las que obtenemos las muestras siguen o no distribuciones Normales, si las varianzas de dichas poblaciones son conocidas o desconocidas, y si la muestras son grandes o no, obtenemos las siguientes expresiones correspondientes a los diferentes intervalos de confianza.

10. Intervalos de Confianza para la Comparación de 2 Poblaciones

Intervalo de confianza para la diferencia de dos medias en poblaciones normales, con varianzas poblacionales conocidas, independientemente del tamaño de la muestra

$$\left(\bar{x}_1 - \bar{x}_2 - z_{\alpha/2} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \bar{x}_1 - \bar{x}_2 + z_{\alpha/2} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

En la figura 10.1 aparece un esquema explicativo de la construcción de este intervalo.



$$P \left(\mu_1 - \mu_2 - z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \leq \bar{x}_1 - \bar{x}_2 \leq \mu_1 - \mu_2 + z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right) = 1 - \alpha$$

Figura 10.1 – Cálculo del intervalo de confianza para la diferencia de medias en poblaciones normales con varianzas conocidas a partir de la distribución de la diferencia de medias muestrales $\bar{x}_1 - \bar{x}_2 \sim N(\mu_1 - \mu_2, \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}})$.

Intervalo de confianza para la diferencia de dos medias en poblaciones normales, con varianzas poblacionales desconocidas, independientemente del tamaño de la muestra

Si aún siendo desconocidas, las varianzas pueden considerarse iguales, el intervalo es:

$$\left(\bar{x}_1 - \bar{x}_2 - t_{\alpha/2}^{n_1+n_2-2} \cdot \sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}, \bar{x}_1 - \bar{x}_2 + t_{\alpha/2}^{n_1+n_2-2} \cdot \sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \right)$$

donde s_p^2 es una cuasivarianza ponderada:

$$s_p^2 = \frac{(n_1 - 1) \cdot \hat{s}_1^2 + (n_2 - 1) \cdot \hat{s}_2^2}{n_1 + n_2 - 2}$$

Si las varianzas, desconocidas, no pueden considerarse como iguales, el intervalo es:

$$\left(\bar{x}_1 - \bar{x}_2 - t_{\alpha/2}^v \cdot \sqrt{\frac{\hat{s}_1^2}{n_1} + \frac{\hat{s}_2^2}{n_2}}, \bar{x}_1 - \bar{x}_2 + t_{\alpha/2}^v \cdot \sqrt{\frac{\hat{s}_1^2}{n_1} + \frac{\hat{s}_2^2}{n_2}} \right)$$

donde ν es el número entero más próximo al valor de la expresión:

$$\frac{\left(\frac{\hat{s}_1^2}{n_1} + \frac{\hat{s}_2^2}{n_2}\right)^2}{\frac{\left(\frac{\hat{s}_1^2}{n_1}\right)^2}{n_1 + 1} + \frac{\left(\frac{\hat{s}_2^2}{n_2}\right)^2}{n_2 + 1}} - 2$$

Si los tamaños muestrales son grandes ($n_1 \geq 30$ y $n_2 \geq 30$) las $t_{\alpha/2}^\nu$ y $t_{\alpha/2}^{n_1+n_2-2}$ pueden sustituirse por $z_{\alpha/2}$.

Intervalo de confianza para la diferencia de dos medias en poblaciones no normales, y muestras grandes ($n_1 \geq 30$ y $n_2 \geq 30$)

En este caso, como ya sucedía con la media muestral, los intervalos para la diferencia de medias son los mismos que sus correspondientes en poblaciones normales y, de nuevo, habría que distinguir si las varianzas son conocidas o desconocidas (iguales o diferentes), lo cual se traduce en que sus correspondientes fórmulas son las mismas que las dadas en los párrafos anteriores. No obstante, por tratarse de muestras grandes, también es válida la aproximación de $t_{\alpha/2}^\nu$ y $t_{\alpha/2}^{n_1+n_2-2}$ por $z_{\alpha/2}$, y habitualmente tan sólo se distingue entre varianzas conocidas y desconocidas.

Para varianzas conocidas:

$$\left(\bar{x}_1 - \bar{x}_2 - z_{\alpha/2} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \bar{x}_1 - \bar{x}_2 + z_{\alpha/2} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right)$$

Y para varianzas desconocidas:

$$\left(\bar{x}_1 - \bar{x}_2 - z_{\alpha/2} \cdot \sqrt{\frac{\hat{s}_1^2}{n_1} + \frac{\hat{s}_2^2}{n_2}}, \bar{x}_1 - \bar{x}_2 + z_{\alpha/2} \cdot \sqrt{\frac{\hat{s}_1^2}{n_1} + \frac{\hat{s}_2^2}{n_2}} \right)$$

Si las poblaciones de partida no son normales y las muestras son pequeñas, no puede aplicarse el Teorema Central de Límite y no se obtienen intervalos de confianza para la diferencia de medias.

Para cualquiera de los anteriores intervalos:

- n_1 y n_2 son los tamaños muestrales.
- \bar{x}_1 y \bar{x}_2 son las medias muestrales.
- σ_1 y σ_2 son las desviaciones típicas poblacionales.
- \hat{s}_1 y \hat{s}_2 son las cuasidesviaciones típicas muestrales: $\hat{s}_1^2 = \frac{\sum (x_{1,i} - \bar{x}_1)^2}{n_1 - 1}$, y análogamente \hat{s}_2^2 .
- $z_{\alpha/2}$ es el valor que deja a su derecha una probabilidad $\alpha/2$ en una distribución Normal tipificada.
- $t_{\alpha/2}^{n_1+n_2-1}$ es el valor que deja a su derecha una probabilidad $\alpha/2$ en una distribución t de Student con $n_1 + n_2 - 1$ grados de libertad.
- $t_{\alpha/2}^\nu$ es el valor que deja a su derecha una probabilidad $\alpha/2$ en una distribución t de Student con ν grados de libertad.

Intervalos de confianza para la media de la diferencia en datos emparejados

En muchas ocasiones hay que estudiar una característica en una población en dos momentos distintos, para estudiar cómo evoluciona con el tiempo, o para analizar la incidencia de algún hecho ocurrido entre dichos momentos.

10. Intervalos de Confianza para la Comparación de 2 Poblaciones

En estos casos se toma una muestra aleatoria de la población y en cada individuo de la misma se observa la característica objeto de estudio en los dos momentos citados. Así se tienen dos conjuntos de datos que no son independientes, pues los datos están emparejados para cada individuo. Por consiguiente, no se pueden aplicar los procedimientos vistos anteriormente, ya que se basan en la independencia de las muestras.

El problema se resuelve tomando para cada individuo la diferencia entre ambas observaciones. Así, la construcción del intervalo de confianza para la diferencia de medias, se reduce a calcular el intervalo de confianza para la media de la variable diferencia. Además, si cada conjunto de observaciones sigue una distribución Normal, su diferencia también seguirá una distribución Normal.

Intervalos de confianza para la diferencia de dos proporciones poblacionales p_1 y p_2

Para muestras grandes ($n_1 \geq 30$ y $n_2 \geq 30$) y valores de p_1 y p_2 (probabilidad de “éxito”) cercanos a 0,5, las correspondientes distribuciones Binomiales pueden aproximarse mediante distribuciones Normales de medias respectivas $n_1 p_1$ y $n_2 p_2$, y desviaciones típicas respectivas $\sqrt{n_1 p_1 (1 - p_1)}$ y $\sqrt{n_2 p_2 (1 - p_2)}$. En la práctica, para que sea válida dicha aproximación, se toma el criterio de que tanto $n_1 p_1$ y $n_2 p_2$ como $n_1 (1 - p_1)$ y $n_2 (1 - p_2)$ deben ser mayores que 5. Lo anterior hace que también podamos construir intervalos de confianza para la diferencia de proporciones tomando éstas como medias de variables dicotómicas en las que la presencia o ausencia de la característica objeto de estudio (“éxito” ó “fracaso”) se expresan mediante un 1 ó un 0 respectivamente.

De este modo, en muestras grandes y con distribuciones Binomiales no excesivamente asimétricas (tanto $n_1 p_1$ y $n_2 p_2$ como $n_1 (1 - p_1)$ y $n_2 (1 - p_2)$ deben ser mayores que 5), si denominamos \hat{p}_1 y \hat{p}_2 a la proporción de individuos que presentan el atributo estudiado en la primera y segunda muestras respectivamente, entonces el intervalo de confianza para la diferencia de proporciones con un nivel de significación α viene dado por:

$$\left(\hat{p}_1 - \hat{p}_2 - z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}_1 \cdot (1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2 \cdot (1 - \hat{p}_2)}{n_2}}, \hat{p}_1 - \hat{p}_2 + z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}_1 \cdot (1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2 \cdot (1 - \hat{p}_2)}{n_2}} \right)$$

donde:

- n_1 y n_2 son los respectivos tamaños muestrales.
- \hat{p}_1 y \hat{p}_2 son las proporciones de individuos que presentan los atributos estudiados en sus respectivas muestras.
- $z_{\alpha/2}$ es el valor que deja a su derecha una probabilidad $\alpha/2$ en una distribución Normal tipificada.

En muestras pequeñas o procedentes de unas distribuciones Binomiales fuertemente asimétricas ($n_1 p_1 \leq 5$, $n_2 p_2 \leq 5$, $n_1 (1 - p_1) \leq 5$ ó $n_2 (1 - p_2) \leq 5$) no puede aplicarse el Teorema Central del Límite y la construcción de intervalos de confianza debe realizarse basándose en la distribución Binomial.

Intervalo de Confianza para la Razón de dos Varianzas σ_1^2 y σ_2^2 de Poblaciones Normales

Como ya hemos visto en la sección de los intervalos de confianza para la diferencia de dos medias en poblaciones normales con varianzas desconocidas, los mismos dependen de si las varianzas, aún siendo desconocidas, pueden considerarse iguales o no. Para dar respuesta a esta cuestión, previa al cálculo del intervalo para la diferencia de medias, se construye un intervalo para la razón (cociente) de varianzas de ambas poblaciones. Para ello tenemos en cuenta que si partimos de dos variables X_1 y X_2 que siguen distribuciones normales con varianzas σ_1^2 y σ_2^2 respectivamente, y tomamos muestras de tamaños n_1 y n_2 de las respectivas poblaciones se tiene que la variable

$$F = \frac{\frac{\hat{S}_1^2}{\sigma_1^2}}{\frac{\hat{S}_2^2}{\sigma_2^2}}$$

sigue una distribución F de Fisher de $n_1 - 1$ grados de libertad en el numerador y $n_2 - 1$ grados de libertad en el denominador.

De lo anterior se deduce que el intervalo de confianza con nivel de significación α para $\frac{\sigma_2^2}{\sigma_1^2}$ es

$$\left(\frac{\hat{s}_2^2}{\hat{s}_1^2} \cdot F_{1-\alpha/2}^{(n_1-1, n_2-1)}, \frac{\hat{s}_2^2}{\hat{s}_1^2} \cdot F_{\alpha/2}^{(n_1-1, n_2-1)} \right)$$

Si dentro del intervalo de confianza obtenido está el número 1 (el cociente de varianzas vale la unidad), no habrá, por tanto, evidencia estadística suficiente, con un nivel de significación α , para rechazar que las varianzas sean iguales.

2 Ejercicios resueltos

1. Para ver si una campaña de publicidad sobre un fármaco ha influido en sus ventas, se tomó una muestra de 8 farmacias y se midió el número de unidades de dicho fármaco vendidas durante un mes, antes y después de la campaña, obteniéndose los siguientes resultados:

Antes	147	163	121	205	132	190	176	147
Después	150	171	132	208	141	184	182	145

- a) Crear un conjunto de datos con las variables antes y despues.
- b) Obtener un resumen estadístico en el que aparezcan la media y la desviación típica de ambas variables. A la vista de los resultados: ¿son las medias diferentes?, ¿ha aumentado la campaña el nivel de ventas?, ¿crees que los resultados son estadísticamente significativos?



- 1) Seleccionar el menú Teaching ▶ Estadística descriptiva ▶ Estadísticos.
- 2) En el cuadro de diálogo que aparece seleccionar las variables antes y despues.
- 3) En la solapa Estadísticos básicos activar la casilla de selección para la Media y la Desviación típica y hacer click en el botón Aceptar.

- c) Obtener los intervalos de confianza para la media de la diferencia entre ambas variables con niveles de significación 0,05 y 0,01.



- 1) Seleccionar el menú Teaching ▶ Test paramétricos ▶ Medias ▶ Test t para dos muestras pareadas.
- 2) En el cuadro de diálogo que aparece seleccionar la variable antes en el campo Comparar, la variable después en el campo Con.
- 3) En la solapa Opciones de contraste introducir 0,95 en el campo Nivel de confianza y hacer click en el botón Aceptar.
- 4) Repetir los pasos para el intervalo de confianza con nivel de significación 0,01 poniendo 0,99 en el campo Nivel de confianza.

- d) ¿Existen pruebas suficientes para afirmar con un 95 % de confianza que la campaña de publicidad ha aumentado las ventas? ¿Y si cambiamos los dos últimos datos de la variable despues y ponemos 190 en lugar de 182 y 165 en lugar de 145? Observar qué le ha sucedido al intervalo para la diferencia de medias y darle una explicación.



- 1) En la ventana de edición de datos, cambiar los datos de las dos últimas farmacias y cerrar la ventana.
 - 2) Repetir los pasos del apartado anterior.
- Existen diferencias entre las medias con el nivel de confianza fijado siempre que el intervalo resultante no contenga el valor 0.

2. Una central de productos lácteos recibe diariamente la leche de dos granjas X e Y. Para analizar la calidad de la leche, durante una temporada, se controla el contenido de materia grasa de la leche

que proviene de ambas granjas, con los siguientes resultados:

X		Y	
0,34	0,34	0,28	0,29
0,32	0,35	0,30	0,32
0,33	0,33	0,32	0,31
0,32	0,32	0,29	0,29
0,33	0,30	0,31	0,32
0,31	0,32	0,29	0,31
		0,33	0,32
		0,32	0,33

- Crear un conjunto de datos con las variables grasa y granja.
- Calcular el intervalo de confianza para el cociente de varianzas del contenido de materia grasa de la leche procedente de ambas granjas.



- Seleccionar el menú Teaching ▶ Test paramétricos ▶ Varianzas ▶ Test F de Fisher.
- En el cuadro de dialogo que aparece seleccionar la variable grasa al campo Comparar y seleccionar la variable granja al campo Según.
- En la solapa Opciones de contraste introducir 0,95 en el campo Nivel de confianza y hacer click sobre el botón Aceptar.

Se mantiene la hipótesis de igualdad de varianzas con la confianza fijada si el intervalo resultante contiene el valor 1.

- Calcular el intervalo de confianza con un 95 % de confianza para la diferencia en el contenido medio de materia grasa de la leche procedente de ambas granjas.



- Seleccionar el menú Teaching ▶ Test paramétricos ▶ Medias ▶ Test t para muestras independientes.
- En el cuadro de dialogo que aparece seleccionar la variable grasa al campo Comparar y seleccionar la variable granja al campo Según.
- En la solapa Opciones de contraste introducir 0,95 en el campo Nivel de confianza, marcar la opción Si en el campo Suponer varianzas iguales y hacer click sobre el botón Aceptar.

- A la vista del intervalo obtenido en el punto anterior, ¿se puede concluir que existen diferencias significativas en el contenido medio de grasa según la procedencia de la leche? Justificar la respuesta.



Existen diferencias entre las medias con el nivel de confianza fijado siempre que el intervalo resultante no contenga el valor 0.

- En una encuesta realizada en una facultad, sobre si el alumnado utiliza habitualmente (al menos una vez a la semana) la biblioteca de la misma, se han obtenido los siguientes resultados:

10. Intervalos de Confianza para la Comparación de 2 Poblaciones

Alumno	1	2	3	4	5	6	7	8	9	10	11	12
Respuesta	No	Si	No	No	No	Si	No	Si	Si	Si	Si	No
Sexo	H	M	M	H	H	H	M	M	M	M	H	H

Alumno	13	14	15	16	17	18	19	20	21	22	23
Respuesta	Si	No	Si	No	No	No	Si	Si	Si	No	No
Sexo	M	H	M	H	H	M	H	M	M	M	H

Alumno	24	25	26	27	28	29	30	31	32	33	34
Respuesta	Si	No	No	Si	Si	No	No	Si	No	Si	No
Sexo	M	H	H	M	M	H	H	M	M	M	H

- Crear un conjunto de datos con las variables respuesta y sexo.
- ¿Existen diferencias significativas entre las proporciones de chicos y chicas que usan habitualmente la biblioteca? Justificar la respuesta.



- Seleccionar el menú Teaching ▶ Test paramétricos ▶ Proporciones ▶ Test para la comparación de dos proporciones.
- En el cuadro de dialogo que aparece seleccionar la variable respuesta al campo Comparar, seleccionar la variable sexo al campo Según, introducir el valor si al campo Categoría y hacer click sobre el botón Aceptar.

Hay diferencias entre las proporciones con el nivel de confianza fijado si el intervalo resultante no contiene el valor 0.

- Un profesor universitario ha tenido dos grupos de clase a lo largo del año: uno con horario de mañana y otro de tarde. En el de mañana, sobre un total de 80 alumnos, han aprobado 55; y en el de tarde, sobre un total de 90 alumnos, han aprobado 32. ¿Existen diferencias significativas en el porcentaje de aprobados en ambos grupos? ¿Pueden ser debidas al turno horario? Justificar la respuesta.



- Seleccionar el menú Teaching ▶ Test paramétricos ▶ Proporciones ▶ Test para la comparación de dos proporciones.
- En el cuadro de diálogo que aparece seleccionar la opción Introducción manual de frecuencias, introducir 55 en el campo Frecuencia muestral 1, introducir 80 en el campo Tamaño muestral 1, introducir 32 en el campo Frecuencia muestral 2, introducir 90 en el campo Tamaño muestral 2 y hacer click en el botón Aceptar.

3 Ejercicios propuestos

- Se ha realizado un estudio para investigar el efecto del ejercicio físico en el nivel de colesterol en la sangre. En el estudio participaron once personas, a las que se les midió el nivel de colesterol (en mg/dl) antes y después de desarrollar un programa de ejercicios. Los resultados obtenidos fueron los siguientes:

Nivel Previo	182	232	191	200	148	249	276	213	241	280	262
Nivel Posterior	198	210	194	220	138	220	219	161	210	213	226

- Hallar el intervalo de confianza del 95 % para la diferencia del nivel medio de colesterol antes y después del ejercicio.
- Hallar el intervalo de confianza del 99 % para la diferencia del nivel medio de colesterol antes y después del ejercicio.

- c) A la vista de los intervalos anteriores, ¿se concluye que el ejercicio físico disminuye el nivel de colesterol?
2. En una encuesta realizada en los dos hospitales de una ciudad se pregunta a los pacientes hospitalizados cuando salen del hospital por si consideran que el trato recibido ha sido correcto. En el primero de ellos se pregunta a 200 pacientes y 140 responden que sí, mientras que en el segundo, se pregunta a 300 pacientes y 180 responden que sí.
- a) Calcular el intervalo de confianza para la diferencia de proporciones de pacientes satisfechos con el trato recibido.
 - b) ¿Hay pruebas significativas para un nivel de significación $\alpha = 0,01$ de que el trato recibido en un hospital es mejor que en el otro?
3. El fichero `nations.txt` contiene información sobre el desarrollo de distintos países (tasa de uso de anticonceptivos (`contraception`), producto interior bruto per cápita (`GDP`), tasa de mortalidad infantil (`infant.mortality`) y tasa de fertilidad (`TFR`)). Se pide:
- a) Importar el fichero `nations.txt` en un conjunto de datos.
 - b) Crear una nueva variable `nivel.economico` que tome el valor `Ricos` para los países con un producto interior bruto per cápita superior a 10000 US\$ y el valor `Pobres` a los países con un producto interior bruto per cápita inferior a dicha cantidad.
 - c) ¿Existen diferencias significativas en el uso de anticonceptivos entre los países ricos y pobres? Justificar la respuesta.
 - d) ¿Existen diferencias significativas en la tasa de fertiliad entre los países ricos y pobres? Justificar la respuesta.
 - e) ¿Existen diferencias significativas en la tasa de mortalidad infantil entre los países ricos y pobres? Justificar la respuesta.