



# INFERENCIA Y MODELOS ESTADÍSTICOS

Jacqueline Köhler C. y José Luis Jara V.



# Tabla de contenido

<b>1</b>	<b>Introducción</b>	<b>1</b>
1.1	Inferencia y modelos estadísticos . . . . .	1
1.2	Variables, parámetros y estadísticos . . . . .	3
1.3	Conociendo R . . . . .	5
1.3.1	Importación de datos . . . . .	5
1.3.2	Importación de paquetes . . . . .	7
1.3.3	Construcción de una matriz de datos . . . . .	7
1.3.4	Modificación de una matriz de datos . . . . .	8
1.3.5	Fórmulas . . . . .	12
1.4	Ejercicios propuestos . . . . .	13
<b>2</b>	<b>Exploración de datos</b>	<b>15</b>
2.1	Estadísticas descriptivas . . . . .	15
2.1.1	Estadísticas descriptivas para datos numéricos . . . . .	16
2.1.2	Estadísticas descriptivas para datos categóricos . . . . .	19
2.1.3	Trabajando con datos agrupados . . . . .	23
2.2	Representación gráfica de datos . . . . .	23
2.2.1	Una variable numérica . . . . .	24
2.2.2	Una variable categórica . . . . .	26
2.2.3	Dos variables numéricas . . . . .	28
2.2.4	Dos variables categóricas . . . . .	30
2.2.5	Una variable numérica y otra categórica . . . . .	32
2.3	Ejercicios propuestos . . . . .	33
<b>3</b>	<b>Variables aleatorias y distribuciones de probabilidad</b>	<b>37</b>
3.1	Variables aleatorias . . . . .	37
3.2	Distribuciones continuas . . . . .	41
3.2.1	Distribución normal . . . . .	42
3.2.2	Distribución Z . . . . .	44
3.2.3	Distribución chi-cuadrado . . . . .	45
3.2.4	Distribución t de Student . . . . .	45
3.2.5	Distribución F . . . . .	47
3.3	Distribuciones discretas . . . . .	48
3.3.1	Distribución de Bernoulli . . . . .	48
3.3.2	Distribución geométrica . . . . .	48
3.3.3	Distribución binomial . . . . .	49
3.3.4	Distribución binomial negativa . . . . .	51
3.3.5	Distribución de Poisson . . . . .	52
3.4	Ejercicios propuestos . . . . .	53
<b>4</b>	<b>Fundamentos para la inferencia</b>	<b>55</b>
4.1	Estimadores puntuales . . . . .	55
4.2	Modelos estadísticos . . . . .	57
4.3	Error estándar . . . . .	58

4.4	Intervalos de confianza . . . . .	59
4.5	Pruebas de hipótesis . . . . .	60
4.5.1	Prueba formal de hipótesis con valores $p$ . . . . .	62
4.5.2	El efecto del nivel de significación . . . . .	66
4.6	Inferencia para otros estimadores . . . . .	66
4.6.1	Estimadores puntuales con distribución cercana a la normal . . . . .	67
4.6.2	Estimadores con otras distribuciones . . . . .	67
4.7	Ejercicios propuestos . . . . .	68



# CAPÍTULO 1. INTRODUCCIÓN

Este libro tiene como propósito acompañarte en el aprendizaje de las primeras nociones de inferencia estadística y de creación de modelos estadísticos. En este primer capítulo comenzaremos por buscar definiciones iniciales para los conceptos de inferencia y modelo, para luego abordar algunas nociones iniciales acerca de los datos empleados en estadística y algunas herramientas para que puedas empezar a usar el entorno de programación R, con el cual trabajaremos a lo largo de todo el texto. Te sugerimos, entonces, que lo instales junto con el entorno de desarrollo integrado RStudio.

## 1.1 INFERENCIA Y MODELOS ESTADÍSTICOS

La Real Academia Española (2014) define **inferencia** como “acción y efecto de inferir”. Esto por sí solo no nos dice mucho, pero si buscamos también la definición de **inferir**, encontraremos que significa “deducir algo o sacarlo como conclusión de otra cosa”. A partir de estas definiciones, y de acuerdo con Devore (2008, p. 5), podemos decir que la **estadística inferencial** es una rama de la estadística que busca obtener una conclusión para un conjunto de individuos o elementos (denominado **población**) a partir de información recolectada de un subconjunto de éste (llamado **muestra**).

Llegar a una definición de **modelo estadístico** puede ser bastante más complejo. Como nos muestra la figura 1.1, ¡un modelo puede ser muchas cosas diferentes! Veamos qué nos dice la Real Academia Española (2014):



Figura 1.1: ejemplos de modelos.

1. Arquetipo o punto de referencia para imitarlo o reproducirlo.
2. En las obras de ingenio y en las acciones morales, ejemplar que por su perfección se debe seguir e imitar.
3. Representación en pequeño de alguna cosa.
4. Esquema teórico, generalmente en forma matemática, de un sistema o de una realidad compleja, como la evolución económica de un país, que se elabora para facilitar su comprensión y el estudio de su comportamiento.
5. Objeto, aparato, construcción, etc., o conjunto de ellos realizados con arreglo a un mismo diseño. *Auto modelo 1976. Lavadora último modelo.*
6. Vestido con características únicas, creado por determinado modista, y, en general, cualquier prenda de vestir que esté de moda.
7. En empresas, u. en aposición para indicar que lo designado por el nombre anterior ha sido creado como ejemplar o se considera que puede serlo. *Empresa modelo. Granjas modelo.*
8. Esc. Figura de barro, yeso o cera, que se ha de reproducir en madera, mármol o metal.
9. *Cuba* impreso (||hoja con espacios en blanco).
10. Persona que se ocupa de exhibir diseños de moda.
11. Persona u objeto que copia el artista.

Puede ser de ayuda tener en cuenta algunas definiciones e ideas que nos ofrece la literatura. Kaplan (2009), por ejemplo, señala que un modelo es una representación con un propósito particular. Pero otros contribuyen a enriquecer esta definición:

- Representación simplificada de la realidad en la que aparecen algunas de sus propiedades (Joly, 1988).
- Permiten estudiar de forma simple y comprensible una porción de la realidad (Ríos, 1995).
- Dejan cosas fuera y pueden llevar a conclusiones equivocadas (Kaplan, 2009).
- Resumen de manera conveniente, a juicio de los sus creadores, los aspectos más relevantes del fenómeno estudiado y sus relaciones (Mendez Ramírez, 1998).
- Están mediados por el diseño, es decir la forma de seleccionar y observar (o manipular) la realidad modelada (Mendez Ramírez, 2012).
- Son pequeños, económicos, seguros y fáciles de transportar, copiar y modificar (Kaplan, 2009).

Si a esta concepción del significado de la palabra modelo le agregamos nuevos conceptos fundamentales que nos ofrecen otros autores, podemos acercarnos un poco más a la idea de **modelo estadístico**:

- Modelo matemático de la regularidad estadística de los posibles resultados de la evolución de un fenómeno aleatorio (Mendez Ramírez, 1998).
- Distribución de probabilidad construida para poder hacer inferencias o tomar decisiones desde datos (Freedman, 2009).
- Descripción simple de un proceso probabilístico que puede haber dado origen a un conjunto de datos observados (McCullagh, 2002).
- Modelo estocástico que contiene parámetros desconocidos que deben ser estimados en base a suposiciones acerca del modelo y los datos (SAS Institute Inc., 2008).

Pero, ¿para qué sirven los modelos estadísticos? Diversos autores nos muestran que tales modelos son muy útiles en diversos contextos:

- Para describir (Kaplan, 2009) o resumir datos (Freedman, 2009).
- Para clasificar (Kaplan, 2009) o predecir (Freedman, 2009; Kaplan, 2009).
- Para anticipar el resultado de intervenciones (Freedman, 2009; Kaplan, 2009).

Ahora que hemos definido nuestros conceptos iniciales, veamos algunas definiciones y herramientas que nos permitan comenzar a explorarlos.

## 1.2 VARIABLES, PARÁMETROS Y ESTADÍSTICOS

Comencemos a partir de un ejemplo para ilustrar algunas ideas iniciales acerca de los datos. La tabla 1.1 muestra las primeras filas de la matriz de datos **ffaa**<sup>1</sup>, que almacena información acerca de miembros activos de las Fuerzas Armadas de Chile. Cada columna de la matriz representa una **variable** o característica, mientras que cada fila corresponde a una **unidad de observación** o instancia. Así, cada fila de la tabla 1.1 almacena datos de una misma persona. El uso de **matrices de datos** para almacenar los datos es muy conveniente, pues nos ayuda a acceder a los datos y modificarlos más fácilmente. Por ejemplo, para agregar una nueva observación, basta con añadir una nueva fila a la matriz. Si queremos eliminar una característica, simplemente borramos la columna correspondiente. Para consultar una característica en particular de una observación, solo necesitamos conocer la fila y la columna correspondientes.

id	género	estatura	escalafón	servicio	antigüedad	rama
1	M	1,77	S	89,91	15	E
2	M	1,97	O	65,14	30	C
3	F	1,65	O	97,03	12	A
4	M	1,82	S	76,29	9	A
5	F	1,73	S	69,46	7	M
6	M	1,78	S	97,67	21	E
7	M	1,87	O	72,09	27	C
8	F	1,91	S	94,53	11	A
⋮	⋮	⋮	⋮	⋮	⋮	⋮
6051	M	1.72	S	86.48	17	E

Tabla 1.1: algunas filas de la matriz de datos **ffaa**.

Antes de comenzar a trabajar con los datos, tenemos que estar seguros de comprender cada uno de sus aspectos. Para ello, las siguientes preguntas pueden ser un excelente punto de partida:

- ¿A qué corresponde cada característica?
- ¿Cuáles son sus unidades de medición?
- ¿Qué valores puede tomar?

La tabla 1.2 muestra la descripción de las características presentes en la matriz de datos de la tabla 1.1. En ella se explica el significado de cada columna de la matriz junto a su rango de valores o a un listado de valores posibles.

Variable	Descripción
id	Identificador de la observación.
género	Género del estudiante (M: masculino, F: femenino).
estatura	Estatura (m).
escalafón	Escalafón al que pertenece (O: oficial, S: suboficial).
servicio	Evaluación de servicio (entre 0 y 100).
antigüedad	Años que lleva en servicio activo.
rama	Rama de las FF.AA. A la que pertenece (E: Ejército, M: Marina, A: Fuerza Aérea, C: Carabineros).

Tabla 1.2: descripción de las variables para el conjunto de datos **ffaa**.

Si estudias la tabla 1.2 con detalle, podemos notar diferencias interesantes entre las variables, más allá de su descripción. Por ejemplo, no todas ellas pueden tomar los mismos valores. Con esto aparece la noción de **tipos de variables**, los cuales podemos jerarquizar:

<sup>1</sup>Los datos aquí presentados son ficticios y han sido creados únicamente con fines pedagógicos.

- **Numéricas:** pueden tomar muchos valores numéricos, y son sensibles a operaciones aritméticas. Pueden separarse en:
  - **Continuas:** pueden tomar cualquier valor (en un intervalo) del conjunto de los reales. Por ejemplo, las variables **estatura** y **servicio** descritas en la tabla 1.2.
  - **Discretas:** no es posible que tomen cualquier valor (en un intervalo). Por ejemplo, podrían tomar únicamente valores enteros no negativos, como la variable **antigüedad** de la matriz de datos **ffaa**.
- **Catóricas:** solo pueden tomar un valor de entre un conjunto acotado. Cada posible valor se denomina **nivel**. Entre las variables catóricas es posible distinguir variables:
  - **Nominales:** no existe un orden natural entre los niveles. Ejemplos de variables nominales son **género** y **rama** de la matriz de datos **ffaa**.
  - **Ordinales:** existe un orden natural entre los niveles. Por ejemplo, la idea de jerarquía es evidente al distinguir entre oficiales y suboficiales en la variable **escalafón** de la tabla 1.1.

Tener diferentes tipos de variables significa que debemos medirlas con distintas clases de escalas, las cuales se distinguen por sus propiedades y los tipos de operaciones que permiten:

- **Escala nominal:** sirve solo para separar un conjunto de elementos en subclases excluyentes entre sí. Los valores no son más que nombres o estados, por lo que no podemos hacer operaciones aritméticas ni podemos establecer relaciones de orden.
- **Escala ordinal o de rangos:** esta escala, al igual que la nominal, permite separar un conjunto de elementos en subclases excluyentes entre sí. Una vez más, los valores son solo nombres o estados, por lo que tampoco podemos hacer operaciones aritméticas. Pero en este caso sí podemos establecer una relación de orden, aunque para ello es necesario que la variable tenga a lo menos tres niveles. A modo de ejemplo, si queremos una variable para medir el nivel de estudios de las personas en un grupo demográfico, podríamos considerar una escala ordinal con los niveles “ninguna”, “básica completa”, “media”, “superior” y “postgrado”. Aquí podemos apreciar claramente que los niveles están ordenados de manera creciente.
- **Escala de intervalo:** sirve para datos continuos o discretos con una gran cantidad de niveles. Además de la noción de orden de la escala ordinal, se cumple que la distancia entre dos valores cualesquiera de la escala es conocida y constante, por lo que podemos emplear operaciones aritméticas. Aunque el punto cero y la unidad de medida son arbitrarios, la razón entre dos intervalos es independiente de ambos elementos. Tomemos, por ejemplo, la escala Celsius de temperatura. El cero está dado por el punto de congelación del agua. La medida o tamaño se calcula en base a los puntos de congelación y ebullición del agua. Sin embargo, a pesar de estos parámetros arbitrarios, el cambio en la cantidad de calor es el mismo si aumentamos la temperatura de 10 a 15 grados Celsius, o de 25 a 30. Si miramos ahora la escala Fahrenheit de temperatura, los puntos fijos son diferentes a los empleados por la escala Celsius, por lo que el cero no significa lo mismo. Sin embargo, existe una transformación lineal que nos permite transformar una medida en una escala a su equivalente en otra escala.
- **Escala de razón:** cumple con todos los atributos de la escala de intervalos, pero además tiene su origen en un cero verdadero. Ejemplos de tales escalas son, por ejemplo, las que permiten medir la masa o la distancia. En una escala de razón, la diferencia entre dos puntos es independiente de la unidad de medida. Por ejemplo, si medimos la masa de dos objetos, la razón es constante independientemente de si empleamos kilogramos, libras u onzas (a diferencia de lo que ocurre con la temperatura usando las escalas Celsius y Fahrenheit).

La estadística usa los datos para responder diversas preguntas, muchas de las cuales se orientan a encontrar relaciones entre variables. Así, dos variables pueden ser:

1. **Independientes:** no existe asociación o relación entre las variables.
2. **Dependientes:** existe una asociación o relación entre las variables. Puede existir:
  - **Asociación positiva:** si una variable crece, la otra también lo hace.
  - **Asociación negativa:** si una variable crece, la otra decrece.

En el contexto de la estadística, decimos que un **parámetro** es cualquier número que describa una población en forma resumida, como por ejemplo la media poblacional. A su vez, un **estadístico** es “cualquier cantidad cuyo valor puede ser calculado a partir de datos muestrales” Devore (2008, p. 204), como por ejemplo la media, la mediana o la desviación estándar de un conjunto de datos observados. Si bien a primera vista



ambos conceptos parecen similares, en realidad existe una diferencia importante entre ellos: el parámetro describe una población, mientras que el estadístico, al ser calculado a partir de una muestra, no es más que una **estimación puntual** del parámetro.

Si necesitas más ejemplos o quieres complementar lo aprendido, puedes consultar los textos de referencia para esta sección. Diez, Barr y Çetinkaya-Rundel (2017, pp. 9-19) describe los principales conceptos relativos a datos, tipos de variables y relaciones entre variables. En Dagnino (2014) puedes aprender más sobre escalas de medición.

## 1.3 CONOCIENDO R

R es un ambiente de software gratuito para estadística computacional y elaboración de gráficos. En esta sección conoceremos algunas herramientas que nos ayudarán a lo largo de este libro. Desde luego, estas breves páginas no pretenden ser un tutorial completo del lenguaje, sino más bien un punto de partida para que podamos aplicar los contenidos que aquí se abordan. Como ya señalamos, sugerimos el uso del entorno integrado de desarrollo RStudio, cuya documentación e instrucciones de instalación podemos consultar en RStudio (2021). En The R Foundation (s.f.) y Carchedi, De Mesmaeker y Vannoorenberghe (s.f.) podemos encontrar documentación acerca del lenguaje R y sus paquetes.

### 1.3.1 Importación de datos

Una de las primeras cosas que necesitamos conocer es cómo importar o cargar una matriz de datos (denominada *data frame* en R) desde un archivo de texto plano (.txt) o de valores separados por coma (.csv). Para lograrlo con éxito, debemos tener en cuenta algunas orientaciones para preparar los datos adecuadamente:

- La primera fila se usa para los nombres de las columnas o variables.
- La primera columna contiene los nombres de las observaciones, que deben ser únicos.
- Los nombres de las columnas deben respetar las convenciones de R:
  - No está permitido el uso de espacios ni símbolos especiales (?, \$, \*, +, #, (, ), -, /, }, {, |, >, <, etc.). Solo se admite el uso de puntos (.) y guiones bajos (\_).
  - Los nombres de variables no pueden comenzar con un dígito.
  - Los nombres de las columnas deben ser únicos.
- R es sensible a las mayúsculas.
- No puede haber filas en blanco.
- No debe tener comentarios.
- Los valores faltantes deben ser denotados mediante NA.
- Para columnas con fechas, se usa el formato mm/dd/aaaa.
- El archivo debe tener uno de los siguientes formatos, ejemplificados en la figura 1.2:
  - Extensión .txt con tabulaciones como delimitador y punto decimal para valores flotantes.
  - Extensión .csv en formato inglés, con comas (,) como delimitador y punto decimal para valores flotantes.
  - Extensión .csv en formato español, con punto y comas (;) como delimitador y coma decimal para valores flotantes.

El script 1.1 muestra las diferentes funciones para importar datos en R, donde las líneas que comienzan por # corresponden a comentarios. La línea 2 carga el conjunto de datos `mtcars`, disponible en R, mientras que las líneas 5, 9 y 16 importan datos desde archivos. Tanto `read.delim()` como `read.csv()` y `read.csv2()` se usan

id	género	estatura	escalafón	servicio	antigüedad	rama
1	M	1.77	S	89.91	15	E
2	M	1.97	O	65.14	30	C
3	F	1.65	O	97.03	12	A
4	M	1.82	S	76.29	9	A
5	F	1.73	S	69.46	7	M
6	M	1.78	S	97.67	21	E
7	M	1.87	O	72.09	27	C
8	F	1.91	S	94.53	11	A

(a) Texto plano delimitado por tabulaciones.

```
id,género,estatura,escalafón,servicio,antigüedad,rama
1,M,1.77,S,89.91,15,E
2,M,1.97,O,65.14,30,C
3,F,1.65,O,97.03,12,A
4,M,1.82,S,76.29,9,A
5,F,1.73,S,69.46,7,M
6,M,1.78,S,97.67,21,E
7,M,1.87,O,72.09,27,C
8,F,1.91,S,94.53,11,A
```

(b) Valores separados por comas (inglés).

```
id;género;estatura;escalafón;servicio;antigüedad;rama
1;M;1.77;S;89.91;15;E
2;M;1.97;O;65.14;30;C
3;F;1.65;O;97.03;12;A
4;M;1.82;S;76.29;9;A
5;F;1.73;S;69.46;7;M
6;M;1.78;S;97.67;21;E
7;M;1.87;O;72.09;27;C
8;F;1.91;S;94.53;11;A
```

(c) Valores separados por punto y comas (español).

Figura 1.2: formatos de archivo para importar datos en R.

de la misma forma, pudiendo recibir como argumento una llamada al selector de archivos (`file.choose()`), como en la línea 5, o la ruta completa para el archivo, como en la línea 9. En el caso de la línea 16, basta con proporcionar el nombre de archivo pues la función `setwd()` (línea 12) permite establecer el directorio de trabajo de R para la sesión. Las funciones `head()` y `tail()` (líneas 20 y 24) proporcionan una buena manera de inspeccionar los datos cargados, pues muestran por consola las primeras y últimas filas de la matriz de datos, respectivamente.

Script 1.1: sentencias para importar un conjunto de datos.

```
1 # Cargar un conjunto de datos disponible en R.
2 datos1 <- mtcars
3
4 # Importar desde un archivo de texto plano delimitado por tabuladores.
5 datos2 <- read.delim(file.choose())
6
7 # Importar desde un archivo de valores separados por coma
8 # en formato inglés (figura 1.2 b).
9 datos3 <- read.csv("C:\\Inferencia\\ejemplo1-csv-eng.csv")
10
11 # Configurar carpeta de trabajo
12 setwd("C:\\Inferencia")
13
14 # Importar desde un archivo de valores separados por coma
15 # en formato español (figura 1.2 c).
16 datos4 <- read.csv2("ejemplo1-csv-esp.csv")
17
18 # Mostrar las primeras 6 filas del conjunto de datos
19 # almacenado en datos1.
20 head(datos1)
21
22 # Mostrar las últimas 6 filas del conjunto de datos
23 # almacenado en datos1.
24 tail(datos1)
```

### 1.3.2 Importación de paquetes

Si bien el entorno R básico incluye muchísimas funcionalidades, existe una enorme variedad de paquetes o colecciones que incorporan otras nuevas o mejoran las ya existentes.

Antes de usar un paquete por primera vez tenemos que instalarlo. Para ello, podemos usar la sentencia que se muestra en la línea 2 del script 1.2. Debemos tener en cuenta que la función `install.packages()` requiere que el nombre del paquete se escriba entre comillas.

Para poder usar un paquete, existen las sentencias `library()` (línea 5 del script 1.2) y `require()` (línea 8), que reciben como argumento el nombre del paquete (sin comillas). Si bien ambas sentencias pueden usarse indistintamente, se diferencian en que `library()` termina la ejecución con un mensaje de error si el paquete no está instalado, mientras que `require()` solo emite una advertencia.

Una forma elegante de evitar errores es verificar si un paquete se encuentra instalado antes de usarlo, para lo que podemos usar una combinación de las sentencias anteriores, como muestran las líneas 11 a 14 del script 1.2. Cabe destacar que la opción `dependencies = TRUE` en la línea 12 asegura que se instalen además aquellos paquetes que son requeridos por el que se desea instalar. Fijémonos que el lenguaje de programación R usa **argumentos con nombre**.

Script 1.2: instalar y cargar paquetes de R.

```
1 # Instalar un paquete.
2 install.packages("ggpubr")
3
4 # Primera forma de importar un paquete.
5 library(ggpubr)
6
7 # Segunda forma de importar un paquete.
8 require(ggplot2)
9
10 # Importar un paquete, instalándolo de ser necesario.
11 if(!require(dplyr)){
12   install.packages("dplyr", dependencies=TRUE)
13   require(dplyr)
14 }
```

### 1.3.3 Construcción de una matriz de datos

Consideremos la idea de construir una matriz de datos que contenga el nombre, la fecha de nacimiento y las calificaciones de los estudiantes en las tres evaluaciones de una asignatura. El script 1.3 crea esta matriz de datos en R con tres observaciones. En las líneas 2 a 4 crea un vector de strings con los nombres de los estudiantes y lo almacena en la variable `nombre`. De manera similar, en la línea 8 crea un vector de fechas. Debemos notar que para ello construye un vector de tres strings con las fechas en formato `aaaa-mm-dd`, el cual es entregado como argumento a la función `as.Date()` para que sean convertidos al formato de fecha. Las líneas 12 a 14 crean tres vectores de flotantes para las calificaciones obtenidas por los estudiantes. Hasta este punto, solo se tienen muchas variables con vectores de largo 3, los cuales deben ser combinados para formar una matriz de datos donde cada vector sea una columna. La función `data.frame()`, en las líneas 18 a 22, realiza esta tarea. Dicha función recibe como argumentos tantos vectores como variables tenga el conjunto de datos, y toma los nombres de las variables que los contienen como nombres de las columnas. Cabe destacar que, en la línea 23, `data.frame()` recibe un argumento adicional, el booleano `stringsAsFactors`, con valor

falso. Esto se debe a que, si no se entrega este parámetro, R asume que su valor por defecto es verdadero, por lo que interpreta el vector de strings como una variable categórica y asigna un valor numérico a cada nivel.

La última línea del script 1.3 permite guardar la matriz de datos en un archivo de valores separados por comas (formato español). La función `write.csv2()` recibe como argumentos el nombre de la variable que contiene la matriz de datos y una cadena de caracteres con el nombre del archivo. El argumento `row.names = FALSE` indica que no deseamos guardar los nombres de las filas. Si queremos guardar nuestra matriz de datos en un archivo separado por comas en formato inglés, podemos hacerlo mediante la función `write.csv()`, que funciona del mismo modo que `write.csv2()`.

Script 1.3: construir un dataframe.

```
1 # Crear un vector de strings y guardarlo en la variable nombre.
2 nombre <- c("Alan Brito Delgado",
3             "Zacarías Labarca del Río",
4             "Elsa Payo Maduro")
5
6 # Crear un vector de fechas y guardarlo en la variable
7 # fecha_nacimiento.
8 fecha_nacimiento <- as.Date(c("2008-1-25", "2006-10-4", "2008-3-27"))
9
10 # Crear tres vectores de reales entre 1.0 y 7.0 y guardarlos
11 # en prueba_i, respectivamente.
12 prueba_1 <- c(5.5, 3.4, 4.5)
13 prueba_2 <- c(3.2, 4.7, 4.1)
14 prueba_3 <- c(4.8, 4.3, 5.1)
15
16 # Construir un data frame a partir de los vectores anteriores y
17 # guardarlo en la variable dataframe.
18 dataframe <- data.frame(nombre,
19                           fecha_nacimiento,
20                           prueba_1,
21                           prueba_2,
22                           prueba_3,
23                           stringsAsFactors = FALSE)
24
25 # Guardar un dataframe en un archivo csv (formato español).
26 write.csv2(dataframe, "C:/Inferencia/Ejemplo.csv", row.names = FALSE)
```

### 1.3.4 Modificación de una matriz de datos

Muchas veces tendremos la necesidad de modificar la matriz de datos. Algunas tareas, como agregar o quitar una columna o un observación pueden hacerse de manera bastante sencilla, como ilustra el script 1.4.

Script 1.4: modificaciones sencillas de una matriz de datos.

```
1 # Leer un dataframe desde archivo csv.
2 datos <- read.csv2("C:/Inferencia/Ejemplo.csv", stringsAsFactors = FALSE)
3
4 # Eliminar del data frame la columna fecha_nacimiento.
5 dataframe$fecha_nacimiento <- NULL
6
7 # Agregar al data frame la columna edad.
8 dataframe$edad <- c(23, 25, 23)
9
```

```

10 # Crear una nueva observación.
11 nueva <- data.frame(nombre="Elba Calao del Río",
12                     prueba_1 = 6.4,
13                     prueba_2 = 2.3,
14                     prueba_3 = 4.6,
15                     edad = 24)
16
17 # Agregar la nueva observación al data frame.
18 dataframe <- rbind(dataframe, nueva)
19
20 # Eliminar las primeras 3 observaciones del data frame.
21 dataframe <- dataframe[-c(1:3),]
22
23 # Guardar el dataframe en un archivo csv .
24 write.csv2(dataframe, "C:/Inferencia/Ejemplo_mod.csv", row.names = FALSE)

```

Sin embargo, también podemos vernos en la necesidad de realizar transformaciones más complejas. El paquete `dplyr` ofrece un conjunto de funciones que simplifica esta tarea:

- `filter()`: selecciona instancias (filas) de acuerdo a su valor.
- `arrange()`: modifica el orden de las filas.
- `select()`: permite seleccionar variables (características) por sus nombres, a la vez que las reordena.
- `mutate()`: permite agregar nuevas variables que se obtienen como funciones de otras ya existentes.

Para mostrar el uso de estas funciones (script 1.5) usaremos el conjunto de datos `iris`, disponible en R. Este contiene 150 observaciones pertenecientes a tres especies de una flor llamada iris: setosa, versicolor y virginica, para las cuales se registran el largo y ancho de sus sépalos y de sus pétalos (en centímetros). Puedes consultar otras funciones y ejemplos más detallados en Müller (2021) y Wickham y Golemund (2017, cap. 5).

Script 1.5: modificación de una matriz de datos con el paquete `dplyr`.

```

1 library(dplyr)
2
3 # Cargar dataframe iris incluido en R.
4 datos <- iris
5
6 # Seleccionar observaciones correspondientes a la especie versicolor.
7 versicolor <- datos %>% filter(Species == "versicolor")
8
9 # Seleccionar observaciones de la especie versicolor cuyos sépalos tengan una
10 # longitud igual o superior a 6 cm.
11 largas <- datos %>% filter(Species == "versicolor" & Sepal.Length >= 6)
12
13 # Seleccionar la especie y variables relativas a los pétalos.
14 petalos <- datos %>% select(Species, starts_with("Petal"))
15
16 # Seleccionar variables de ancho y la especie.
17 anchos <- datos %>% select(ends_with("Width"), Species)
18
19 # Agregar al conjunto de datos de los pétalos una nueva variable con la razón
20 # entre el largo y el ancho de éstos.
21 petalos <- petalos %>% mutate(Species, Petal.Width,
22                             Petal.Ratio = Petal.Length / Petal.Width)
23
24 # Ordenar el conjunto de datos de pétalos en forma descendente según la razón
25 # de los pétalos.
26 petalos <- petalos %>% arrange(desc(Petal.Ratio))
27
28 # Ordenar el conjunto de datos de pétalos en forma ascendente según el largo de

```

```

29 # los pétalos.
30 petalos <- petalos %>% arrange(Petal.Length)

```

En el script 1.5 aparece frecuentemente el operador `%>%`, llamado *pipe* y definido en el paquete `magrittr`, cuya función es entregar un valor o el resultado de una expresión a la siguiente llamada a una función. En términos sencillos, la expresión `x%>% f` es equivalente a `f(x)`, y su utilidad es que simplifica la lectura de llamadas a funciones anidadas (Bache, 2014).

Otra transformación que se usa a menudo es la de pivotar la matriz de datos, cuyo efecto es el de “alargar” o “ensanchar” la matriz. En el primer caso, se incrementa la cantidad de filas (observaciones) a la vez que se reduce la cantidad de columnas (variables). Para ello se usa la función `pivot_longer(cols, names_to, values_to)` del paquete `tidyr`, donde:

- `cols`: nombres de las columnas a pivotar.
- `names_to`: especifica el nombre de una nueva columna cuyos valores corresponden a los nombres de las columnas a pivotar.
- `values_to`: especifica el nombre de una nueva columna donde se almacenan los valores de las columnas a pivotar.

En el segundo caso se obtiene como resultado una reducción de la cantidad de filas junto al aumento de la cantidad de columnas. Para ello se usa la función `pivot_wider(names_from, values_from)`, también del paquete `tidyr`, donde:

- `names_from`: especifica el nombre de una variable desde la que se obtienen los nombres de las nuevas columnas.
- `values_from`: especifica el nombre de una variable desde donde se obtienen los valores de las nuevas columnas.

Veamos con un ejemplo el efecto de estas dos transformaciones. El script 1.6 comienza por crear una matriz de datos en que se registran los tiempos de ejecución (en milisegundos) para seis instancias de un problema con cuatro algoritmos diferentes. Las columnas de la matriz de datos original corresponden al identificador de la instancia y cada uno de los algoritmos. Así, la matriz de datos original tiene 6 filas y 5 columnas.

A continuación, se crea una nueva matriz de datos, `datos_largos`, que resulta de pivotar la original para “alargarla”. Al ejecutar el script 1.6 podemos ver que nuestra nueva matriz de datos tiene solo tres columnas, pero que su cantidad de filas es 24. Si miramos con atención, veremos que ahora tenemos 4 filas por cada instancia, una por cada algoritmo (señalado en la columna `Algoritmo`) con su correspondiente tiempo de ejecución (columna `Tiempo`).

Por último, el script 1.6 crea otro conjunto de datos, `datos_anchos`, a partir de `datos_largos`. Al examinar este nuevo conjunto, se puede apreciar que es idéntico al creado inicialmente.

Script 1.6: modificación de una matriz de datos con el paquete `tidyr`.

```

1 library(dplyr)
2 library(tidyr)
3
4 # Crear el data frame.
5 Instancia <- 1:6
6 Quicksort <- c(23.2, 22.6, 23.4, 23.3, 21.8, 23.9)
7 Bubblesort <- c(31.6, 29.3, 30.7, 30.8, 29.8, 30.3)
8 Radixsort <- c(30.1, 28.4, 28.7, 28.3, 29.9, 29.1)
9 Mergesort <- c(25.0, 25.7, 25.7, 23.7, 25.5, 24.7)
10 datos <- data.frame(Instancia, Quicksort, Bubblesort, Radixsort, Mergesort)
11
12 # Mostrar las primeras filas de la matriz de datos.
13 cat("Datos originales\n")
14 print(head(datos))
15 cat("\n")
16

```

```

17 # Convertir la matriz de datos a formato largo.
18 datos_largos <- datos %>% pivot_longer(c("Quicksort", "Bubblesort",
19                                         "Radixsort", "Mergesort"),
20                                         names_to = "Algoritmo",
21                                         values_to = "Tiempo")
22
23 # Mostrar las primeras filas de la matriz de datos largos.
24 cat("Datos largos\n")
25 print(head(datos_largos))
26 cat("\n")
27
28 # Convertir la matriz de datos largos a formato ancho.
29 datos anchos <- datos_largos %>% pivot_wider(names_from = "Algoritmo",
30                                              values_from = "Tiempo")
31
32 # Mostrar las primeras filas de la matriz de datos anchos.
33 cat("Datos anchos\n")
34 print(head(datos anchos))
35 cat("\n")

```

Habrás notado que para poder usar las funciones de `tidyr` se requiere también el paquete `dplyr`. Una alternativa es cargar únicamente el paquete `tidyverse`, el cual los incluye a ambos (entre otros).

Puedes encontrar descripciones más extensas acerca del uso de las funciones del paquete `tidyverse`, junto con ejemplos más avanzados, en Wickham (2021).

En ocasiones puede ser necesario renombrar las columnas para que nos resulte más fácil comprender a qué variable corresponde. La función `rename()` del paquete `dplyr` nos permite hacer esta operación bastante sencilla. Sus argumentos son una lista de elementos de la forma `nuevo nombre = nombre original`. También podemos cambiar el tipo de una variable. Una conversión que nos será muy útil es de variable numérica a categórica, lo que se logra mediante la función `factor(x, levels, labels, ordered)`, donde:

- **x**: nombre de la variable a convertir.
- **levels**: argumento opcional con los posibles valores de la variable categórica.
- **labels**: argumento opcional con las etiquetas asociadas a cada valor.
- **ordered**: valor lógico que especifica si la variable es o no ordinal (falso por defecto).

Tomemos el conjunto de datos `mtcars` (incluido en R) para ejemplificar el uso de estas funciones. La tabla 1.3 muestra la descripción de estos datos. El script 1.7 modifica los nombres de las columnas para que sean más representativos y da formato de variable categórica a las variables que así lo requieren, asignando etiquetas adecuadas para cada nivel.

Variable	Descripción
mpg	Rendimiento, en millas / galón (EEUU).
cyl	Número de cilindros.
disp	Desplazamiento, en pulgadas cúbicas.
hp	Potencia, en caballos de fuerza brutos.
drat	Razón del eje trasero.
wt	Peso, en miles de libras.
qsec	Tiempo que tarda en recorrer un cuarto de milla partiendo desde el reposo, en segundos.
vs	Tipo de motor (0: en forma de V, 1: recto).
am	Tipo de transmisión (0: automática, 1: manual).
gear	Número de marchas hacia adelante.
carb	Número de carburadores.

Tabla 1.3: descripción de las variables para el conjunto de datos `mtcars`.

Script 1.7: modificación del conjunto de datos mtcars para facilitar su comprensión.

```
1 library(dplyr)
2
3 # Cargar conjunto de datos.
4 datos <- mtcars
5
6 # Renombrar columnas.
7 datos <- datos %>% rename(Rendimiento = mpg, Cilindrada = cyl,
8                           Desplazamiento = disp, Potencia = hp,
9                           Eje = drat, Peso = wt, Cuarto_milla = qsec,
10                          Motor = vs, Transmision = am, Cambios = gear,
11                          Carburadores = carb)
12
13 # Dar formato categórico a las variables Motor y Transmision, renombrando
14 # sus niveles.
15 datos[["Motor"]] <- factor(datos[["Motor"]], levels = c(0, 1),
16                           labels = c("V", "Recto"))
17
18 datos[["Transmision"]] <- factor(datos[["Transmision"]], levels = c(0, 1),
19                                 labels = c("Automático", "Manual"))
20
21 # Dar formato ordinal a las variables Cilindrada y Cambios, renombrando
22 # sus niveles.
23 datos[["Cilindrada"]] <- factor(datos[["Cilindrada"]], levels = c(4, 6, 8),
24                                 labels = c("4 cilindros", "6 cilindros",
25                                             "8 cilindros"),
26                                 ordered = TRUE)
27
28 datos[["Cambios"]] <- factor(datos[["Cambios"]], levels = c(3, 4, 5),
29                              labels = c("3 cambios", "4 cambios", "5 cambios"),
30                              ordered = TRUE)
31
32 write.csv2(datos, "C:/Inferencia/Mtcars.csv")
```

### 1.3.5 Fórmulas

Si bien hasta ahora solo tenemos una definición preliminar de lo que es un modelo estadístico, necesitamos conocer una herramienta para representarlos en R, pues son una parte fundamental del funcionamiento de este lenguaje.

Para entender de manera sencilla qué es una fórmula, podemos simplemente decir que permite capturar una expresión no evaluada, y que está asociada a un ambiente. Su sintaxis básica tiene la forma **variable independiente** ~ **variables dependientes**, lo que nos indica, entonces, que las fórmulas representan una relación entre variables.

Tomemos una vez más el conjunto de datos **iris**. Podríamos representar la asociación entre la especie de iris (variable independiente) y las dimensiones de sus pétalos (variables dependientes) como **Species** ~ **Petal.Length** + **Petal.Width**.

Extenderemos las nociones acerca del uso de fórmulas a medida que avancemos en nuestro aprendizaje, pero si quieres aprender más puedes consultar Willems (2017).



## 1.4 EJERCICIOS PROPUESTOS

1. Una encuesta reciente preguntó: “después de la jornada laboral usual, ¿cuántas horas dedica a relajarse o a realizar actividades que disfruta?” a una muestra de 580 chilenas y 575 chilenos. Se encontró que el número promedio de horas era de  $1,30 \pm 0,30$  y  $1,95 \pm 0,25$  para cada grupo, respectivamente.
  - a) ¿Cómo sería una matriz de datos para este estudio? Muestra algunas filas de ella como ejemplos.
  - b) ¿Cuál podría ser la población objetivo?
  - c) ¿Qué se entendería por unidad de observación?
  - d) ¿Qué tipo de variable sería “el número de horas dedicadas a distraerse después de la jornada laboral usual” que respondió cada persona entrevistada?
  - e) ¿Existe alguna variable categórica? Si es así, ¿de qué tipo? ¿Con qué niveles?
  - f) ¿Qué dato(s) correspondería(n) a un estadístico?
  - g) ¿Cuál(es) sería(n) el(los) parámetro(s) en estudio?
  - h) ¿Logra el estudio establecer que ser mujer chilena ocasiona tener menos horas dedicadas a distraerse después de la jornada laboral usual?
2. Investiga para qué sirven y cómo se usan los argumentos `row.names` y `col.names` en las funciones para importar datos desde archivos y la función `data.frame()`.
3. Construye en R una matriz de datos para almacenar las características de una muestra de servidores. Considera a lo menos una variable categórica y una variable numérica.
4. Investiga qué función (o funciones) ofrece R para guardar una matriz de datos en un archivo y úsala(s) para guardar la matriz de datos del ejercicio anterior.
5. Resuelve en R los siguientes ejercicios. Considera para ello el conjunto de datos nativo de R `chickwts`.
  - a) ¿Cómo se puede cargar el conjunto de datos en la variable `pollos`?
  - b) ¿Cómo se ve la estructura de la matriz de datos almacenada en `pollos`?
6. Muestra ejemplos de las distintas transformaciones que se pueden hacer a una matriz de datos usando para ello conjunto de datos nativo de R `ChickWeight`.



## CAPÍTULO 2. EXPLORACIÓN DE DATOS

Siempre es bueno que nos familiaricemos con los datos y algunas de sus características antes de empezar a trabajar con ellos. Esto nos ayuda a decidir qué herramientas son las más adecuadas para dar respuesta a las preguntas que queramos responder. En este capítulo revisaremos las principales estadísticas descriptivas que nos ayudarán a resumir los datos para entenderlos mejor, así como diversos tipos de gráficos que nos permitirán representar los datos de modo que podamos comprenderlos de forma visual. Para ello, tomamos como base los conceptos expuestos en Diez y col. (2017, pp. 26-50), Field y col. (2012, pp. 19-27) y STDHA (s.f.), fuentes que puedes consultar si deseas saber más acerca de estos temas.

Para muchos de los ejemplos de este capítulo usaremos el conjunto de datos `mtcars` con las modificaciones realizadas en el script 1.7, cuyo diccionario de datos se muestra en la tabla 2.1.

Variable	Descripción
Rendimiento	Rendimiento, en millas / galón (EEUU).
Cilindrada	Número de cilindros (4 cilindros, 6 cilindros, 8 cilindros).
Desplazamiento	Desplazamiento, en pulgadas cúbicas.
Potencia	Potencia, en caballos de fuerza brutos.
Eje	Razón del eje trasero.
Peso	Peso, en miles de libras.
Cuarto_milla	Tiempo que tarda en recorrer un cuarto de milla partiendo desde el reposo, en segundos.
Motor	Tipo de motor (V, Recto).
Transmision	Tipo de transmisión (Automático, Manual).
Cambios	Número de cambios hacia adelante (3 cambios, 4 cambios, 5 cambios).
Carburadores	Número de carburadores.

Tabla 2.1: descripción de las variables para el conjunto de datos `mtcars`.

### 2.1 ESTADÍSTICAS DESCRIPTIVAS

Las estadísticas descriptivas son medidas que nos permiten sintetizar y, como su nombre lo indica, describir los datos. Estas pueden aplicarse tanto a una muestra como a una población. Cuando una de estas medidas se aplica a la muestra, corresponde a un **estimador puntual** de la misma medida para la población. Al ser una estimación, no es exacta, aunque la precisión tiende a aumentar mientras mayor sea el tamaño de la muestra.

Un concepto importante a tener en cuenta es la noción de **distribución**. En este capítulo se considera la **distribución de frecuencia**, que representa cuántas veces aparece cada valor para una variable en un conjunto de datos.

### 2.1.1 Estadísticas descriptivas para datos numéricos

Una de las estadísticas descriptivas más empleadas es la **media**, conocida en otros contextos como media aritmética o promedio. Denotamos la **media muestral** por  $\bar{x}$ , donde  $x$  corresponde al nombre de la variable, mientras que para la **media poblacional** empleamos la notación  $\mu_x$ . Esta medida se calcula como muestra la ecuación 2.1, donde  $x_i$  son los  $n$  valores observados de la variable. Podemos entender la media como el punto de equilibrio de la distribución (Diez y col., 2017, p. 28). Así, la media corresponde a una **medida de tendencia central**.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (2.1)$$

El script 2.1 muestra cómo usar la función `mean()` de R para calcular la media de diversas variables del conjunto de datos `mtcars`<sup>1</sup>. Como primer ejemplo, se calcula la media de la variable `Rendimiento`. A continuación se muestra cómo realizar esta operación para dos variables, señaladas por el índice de sus respectivas columnas. Luego, de manera similar, se calculan las medias para cuatro columnas consecutivas de la matriz de datos. En estos dos casos hacemos uso de la función `sapply()`, que permite aplicar una misma función (cualquiera) para múltiples columnas.

Script 2.1: uso de las funciones `mean()` y `sapply()`.

```
1 # Cargar conjunto de datos.
2 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
3                   row.names = 1)
4
5 # Calcular la media para la variable Rendimiento.
6 media <- mean(datos[["Rendimiento"]])
7 cat("Rendimiento medio:", media, "\n\n")
8
9 # Calcular la media para la tercera y quinta columnas
10 # (variables Desplazamiento y Eje).
11 cat("Medias\n")
12 print(sapply(datos[c(3, 5)], mean))
13 cat("\n")
14
15 # Calcular la media para las columnas 3 a 6
16 # (variables Desplazamiento, Potencia, Eje y Peso).
17 cat("Medias\n")
18 print(sapply(datos[3:6], mean))
19 cat("\n")
20
21 # Calcular la media para la variable Rendimiento omitiendo valores faltantes.
22 print(mean(datos[["Rendimiento"]], na.rm = TRUE))
```

La función `mean()` devuelve NA (*not available*, es decir, no disponible) si existen valores faltantes en los datos de entrada. Para prevenir este error, se puede proporcionar un argumento adicional que descarte los valores faltantes, como muestra la última línea del script 2.1.

Una medida de tendencia central alternativa a la media es la **mediana**, que es, simplemente, el valor central de los valores previamente ordenados. Cuando no existe un valor central, vale decir, cuando el tamaño de la muestra es par, la mediana está dada por el promedio simple de los dos valores centrales. En R, la mediana se calcula con la función `median()`.

---

<sup>1</sup>Todas las demás funciones de R mencionadas en esta sección para las que no se proporcione un script se usan del mismo modo que `mean()`.

La **moda** es, simplemente, el valor más frecuente en el conjunto de datos. No obstante, tiene el problema de que puede haber múltiples modas. Dependiendo de la cantidad de modas, se habla de distribuciones **unimodales**, **bimodales** y **multimodales**.

Si bien R no cuenta con una función nativa para encontrar la moda, el paquete **modeest** ofrece la función **mfv()** que entrega el valor más frecuente de una variable. En caso de que dos (o más) valores sean los más frecuentes con igual cantidad de observaciones, los entrega todos en forma de vector.

Las medidas que hemos estudiado hasta ahora buscan describir el centro del conjunto de datos. No obstante, también es importante conocer su **variabilidad o dispersión**, pues así se puede saber qué tan semejantes (o diferentes) son las observaciones entre sí. Estas suelen calcularse en base a la **desviación** de las observaciones, que se entiende como la distancia entre una observación y la media del conjunto de datos. Las dos principales medidas de dispersión son la **varianza** y la **desviación estándar**, ambas basadas en los cuadrados de las distancias, ya que, por una parte, los valores grandes se incrementan más significativamente y, por otra, se opera solo con valores positivos, pues la dirección de la distancia no es de interés.

La varianza muestral se calcula como muestra la ecuación 2.2, donde  $x_i$  son los valores de cada una de las  $n$  observaciones. Cabe destacar que puede emplearse un subíndice para indicar el nombre de la variable, al igual que en el caso de la media.

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (2.2)$$

La desviación estándar de la muestra se define como la raíz cuadrada de la varianza (2.3), medida que resulta de gran utilidad cuando se necesita saber cuán cercanos son los datos a la media, ya que se encuentra en la misma escala que la variable.

$$s = \sqrt{s^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \quad (2.3)$$

Al igual que en el caso de la media, podemos usar las fórmulas anteriores para obtener estimaciones puntuales de la varianza y la desviación estándar de la población, denotadas por  $\sigma^2$  y  $\sigma$ , respectivamente.

Es importante considerar que, si bien la media y la desviación estándar permiten conocer el centro y la dispersión del conjunto de datos, respectivamente, la distribución de los puntos puede ser muy diferente, como ilustra la figura 2.1.

Las funciones de R para calcular la varianza y la desviación estándar son, respectivamente, **var()** y **sd()**.

Aunque menos empleado, el **rango** muestra los valores extremos, es decir, el mínimo y el máximo, de una variable. R ofrece la función **range()** para obtener ambos valores, además de **min()** y **max()** para obtenerlos por separado.

En párrafos anteriores vimos que la mediana es el valor central (o el promedio de los dos valores centrales) del conjunto de datos ordenado, ya sea una población o una muestra. Esto significa, entonces, que esta medida divide el conjunto de datos en dos mitades con igual cantidad de elementos. De manera similar, es posible dividir el conjunto de datos en segmentos más pequeños, por ejemplo en 4, 10 o 100 partes con igual cantidad de elementos. Cada fragmento del conjunto de datos dividido de esta forma recibe el nombre de **cuantil**. Algunas subdivisiones de uso frecuente reciben nombres especiales:

- **Percentiles:** dividen el conjunto de datos en 100 subconjuntos de igual tamaño.
- **Deciles:** dividen el conjunto de datos en 10 subconjuntos de igual tamaño.
- **Quintiles:** dividen el conjunto de datos en 5 subconjuntos de igual tamaño.
- **Cuartiles:** dividen el conjunto de datos en 4 subconjuntos de igual tamaño.

Los cuantiles (al igual que las otras subdivisiones antes mencionadas) se nombran de forma ascendente según el sentido de crecimiento del conjunto de datos. Así, el percentil 1 contiene a los valores más pequeños,

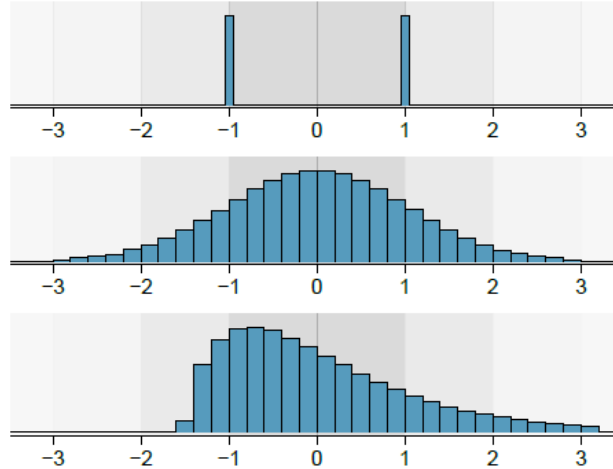


Figura 2.1: tres distribuciones de población muy distintas con media  $\mu = 0$  y desviación estándar  $\sigma = 1$ .  
Fuente: Díez y col. (2017, p. 34).

mientras que el percentil 100, a los más grandes. Cabe destacar que la mediana corresponde al percentil 50 o al cuartil 2, y que nombrar al decil 3 es equivalente al percentil 30.

R proporciona la función `quantile()` para calcular cuantiles, que por defecto calcula los cuartiles, aunque su uso puede generalizarse mediante el parámetro adicional `probs`, como muestra el script 2.2. La función `seq()` genera una secuencia de números equiespaciados, y recibe como argumentos el inicio, el término y el incremento de la secuencia.

Script 2.2: cálculo de cuantiles con la función `quantile()`.

```
1 # Cargar conjunto de datos.
2 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
3                   row.names = 1)
4
5 # Cálculo de percentiles para la variable Rendimiento.
6 cat("Cuartiles:\n")
7 print(quantile(datos[["Rendimiento"]]))
8 cat("\n")
9
10 cat("Quintiles:\n")
11 print(quantile(datos[["Rendimiento"]], seq(0, 1, 0.2)))
12 cat("\n")
13
14 cat("Deciles:\n")
15 print(quantile(datos[["Rendimiento"]], seq(0, 1, 0.1)))
16 cat("\n")
17
18 cat("Percentiles:\n")
19 print(quantile(datos[["Rendimiento"]], seq(0, 1, 0.01)))
```

Ahora que conocemos los cuartiles, podemos introducir una nueva medida de variabilidad que usaremos a menudo, llamada **rango intercuartil** o IQR (por su sigla en inglés), dada por la ecuación 2.4, donde  $Q_1$  y  $Q_3$  corresponden a los cuartiles 1 y 3, respectivamente. Al igual que la varianza y la desviación estándar, mientras más disperso sea el conjunto de datos, mayor será el valor del IQR. En R, la función que calcula este estimador es `IQR()`.

$$IQR = Q_3 - Q_1 \quad (2.4)$$

Muchas veces los conjuntos de datos contienen lo que se conoce como **valores atípicos** o *outliers*. Estos corresponden a observaciones que parecen estar fuera de rango o ser muy extremos con respecto al resto de los datos. Medidas como la media o la desviación estándar son muy sensibles a los valores atípicos, por lo que son propensas a errores ante la presencia de este tipo de observaciones. Para reducir el efecto de los valores extremos muchas veces necesitaremos medidas **robustas**, que son aquellas que proporcionan una estimación confiable aún ante la presencia de valores atípicos. En este escenario, la mediana resulta ser una buena medida de tendencia central y el IQR, una buena medida de dispersión.

Nos encontraremos frecuentemente con la necesidad de calcular varias medidas de tendencia central y de dispersión descritas en el apartado anterior. Por esta razón, R, y algunos de sus paquetes, ofrecen algunas funciones que calculan varios de estos estadísticos con una sola llamada. Tal es el caso de la función nativa `summary()`, que entrega la media, la mediana, el primer y el tercer cuartil, el mínimo y el máximo. Otra función que nos puede ser de mucha ayuda es `summarise()`, del paquete `dplyr`. Con ella podemos calcular varias de las medidas en una sola llamada, como muestra el script 2.3.

Script 2.3: uso de la función `summarise()` del paquete `dplyr`.

```
1 library(dplyr)
2
3 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
4                   row.names = 1)
5
6 # Cálculo de varias medidas para la variable Potencia.
7 medidas_potencia <- datos %>% summarise(Media = mean(Potencia),
8                                         Mediana = median(Potencia),
9                                         Varianza = var(Potencia),
10                                        IQR = IQR(Potencia))
11
12 print(medidas_potencia)
13 cat("\n")
14
15 # Cálculo de la media y la desviación estándar para las variables Peso y
16 # Cuarto_milla.
17 medidas_varias <- datos %>% summarise(Media_P = mean(Peso),
18                                       Media_C = median(Cuarto_milla),
19                                       SD_P = sd(Peso),
20                                       SD_C = sd(Cuarto_milla))
21
22 print(medidas_varias)
23 cat("\n")
```

### 2.1.2 Estadísticas descriptivas para datos categóricos

Cuando queremos trabajar con datos categóricos, medidas como la media o la desviación estándar carecen de sentido. En consecuencia, necesitamos otros estadísticos para resumir el conjunto de datos.

Como primer estadístico para variables categóricas podemos mencionar la **frecuencia**, que corresponde a la cantidad de veces que podemos encontrar cada nivel de la variable en los datos. Otro estadístico importante corresponde a la **proporción**, que corresponde a la frecuencia relativa. En otras palabras, la proporción corresponde a frecuencia de un nivel de la variable dividida por la cantidad total de observaciones.

La mejor alternativa para este tipo de datos es la **tabla de contingencia**, también llamada **matriz de confusión** o **tabla de frecuencias**, donde cada fila representa la cantidad de veces en que ocurre una combinación de variables. También es posible usar porcentajes o proporciones en lugar de la cantidad de

ocurrencia, en cuyo caso se habla de una **tabla de frecuencias relativas**. La tabla 2.2 muestra la tabla de contingencia (de frecuencias) para la variable **Cambios**. Se puede observar, por ejemplo, que el conjunto de datos contiene una muestra de 32 automóviles y que 15 de ellos tienen tres cambios.

3 cambios	4 cambios	5 cambios	Total
15	12	5	32

Tabla 2.2: tabla de contingencia para la cantidad de cambios de los automóviles.

Desde luego, podemos construir tablas de contingencia de manera bastante sencilla en R. El script 2.4 muestra dos formas de obtener la tabla 2.2. La primera es la función `table()` y la segunda, la función `xtabs()`. El funcionamiento de ambas es equivalente, aunque `xtabs()` muestra el nombre de la variable tabulada al imprimir los resultados y `table()` no lo hace. Las tablas entregadas por estas funciones no incluyen los totales por filas, pero la función `marginSums()` permite calcularlos y mostrarlos como un vector. A su vez, la función `addmargins()` permite calcular dichos totales e incorporarlos a la tabla. Para terminar, el las últimas sentencias del script 2.4 ilustran la manera de obtener las tablas de frecuencias relativas con proporciones y porcentajes, respectivamente.

Podemos ver que las llamadas a `table()` y a `xtabs()` son algo diferentes. La primera recibe como argumento la columna de la matriz de datos, es decir, un vector con los datos a tabular, mientras que la segunda recibe una fórmula en que no existe una variable dependiente y la variable categórica es la independiente.

Script 2.4: tabla de contingencia para la variable **Cambios**.

```

1 # Cargar datos.
2 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
3                   row.names = 1)
4
5 # Crear tabla de contingencia para la variable gear.
6 contingencia <- table(datos[["Cambios"]])
7 cat("Tabla de contingencia generada con table():\n")
8 print(contingencia)
9 cat("\n")
10
11 # Otra forma de crear la misma tabla.
12 contingencia <- xtabs(~ Cambios, data = datos)
13 cat("Tabla de contingencia generada con xtabs():\n")
14 print(contingencia)
15 cat("\n")
16
17 # Calcular totales por fila y mostrarlos por separado.
18 totales <- marginSums(contingencia)
19 cat("Totales por fila:\n")
20 print(totales)
21 cat("\n")
22
23 # Calcular totales por fila y agregarlos a la tabla.
24 con_totales <- addmargins(contingencia, 1)
25 cat("Tabla de contingencia con totales por fila:\n")
26 print(con_totales)
27 cat("\n")
28
29 # Convertir a tabla de proporciones
30 proporciones <- prop.table(contingencia)
31 proporciones <- addmargins(proporciones, 1)
32 cat("Tabla de contingencia con proporciones:\n")
33 print(proporciones)
34 cat("\n")
35

```



```

36 # Convertir a tabla de porcentajes con 2 decimales.
37 porcentajes <- round(prop.table(contingencia), 4) * 100
38 porcentajes <- addmargins(porcentajes)
39 cat("Tabla de contingencia con porcentajes:\n")
40 print(porcentajes)
41 cat("\n")

```

También podemos construir matrices de confusión para dos variables categóricas, como muestra la tabla 2.3 para las variables Cambios y Transmisión.

		Cambios			Total
		3 cambios	4 cambios	5 cambios	
Transmision	Automático	15	4	0	19
	Manual	0	8	5	13
	Total	15	12	5	32

Tabla 2.3: tabla de contingencia para las variables Cambios y Transmisión.

En ocasiones resulta útil determinar las proporciones por fila o por columna, que podemos obtener dividiendo el valor de una celda de la matriz por el total de su fila o columna, según corresponda. Así, el total de cada fila (o columna) es igual a 1. Puesto que las proporciones por fila y por columna no son equivalentes, debemos ser cuidadosos al escoger la más adecuada en cada caso. Las tablas 2.4 a 2.6 muestran las proporciones por fila, por columna y generales para la matriz de confusión de la tabla 2.3. La construcción en R de la tabla de contingencia y las tablas de proporciones para dos variables se muestra en el script 2.5.

		Cambios			Total
		3 cambios	4 cambios	5 cambios	
Transmision	Automático	0,7894737	0,2105263	0,0000000	1,0000000
	Manual	0,0000000	0,6153846	0,3846154	1,0000000

Tabla 2.4: tabla de proporciones con totales por fila para la tabla 2.3.

		Cambios		
		3 cambios	4 cambios	5 cambios
Transmision	Automático	1,0000000	0,3333333	0,0000000
	Manual	0,0000000	0,6666667	1,0000000
	Total	1,0000000	0,0000000	1,0000000

Tabla 2.5: tabla de proporciones con totales por columna para la tabla 2.3.

Script 2.5: tablas de contingencia y proporciones para dos variables.

```

1 # Cargar datos.
2 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
3                   row.names = 1)
4
5 # Crear tabla de contingencia para las variables Transmision y gear.
6 contingencia <- table(datos[["Transmision"]], datos[["Cambios"]])
7 cat("Tabla de contingencia generada con table():\n")
8 print(contingencia)
9 cat("\n")
10
11 # Otra forma de crear la misma tabla.
12 contingencia <- xtabs(~ Transmision + Cambios, data = datos)
13 cat("Tabla de contingencia generada con xtabs():\n")
14 print(contingencia)
15 cat("\n")

```

		Cambios			Total
		3 cambios	4 cambios	5 cambios	
Transmision	Automático	0,46875	0,12500	0,00000	0,59375
	Manual	0,00000	0,25000	0,15625	0,40625
	Total	0,46875	0,37500	0,15625	1,00000

Tabla 2.6: tabla de proporciones con totales por fila y columna para la tabla 2.3.

```

16
17 # Proporciones con totales por fila.
18 proporciones_fila <- prop.table(contingencia, margin=1)
19 proporciones_fila <- addmargins(proporciones_fila, margin=2)
20 cat("Tabla de contingencia con proporciones totales por fila:\n")
21 print(proporciones_fila)
22 cat("\n")
23
24 # Proporciones con totales por columna.
25 proporciones_columna <- prop.table(contingencia, margin=2)
26 proporciones_columna <- addmargins(proporciones_columna, margin=1)
27 cat("Tabla de contingencia con proporciones totales por columna:\n")
28 print(proporciones_columna)
29 cat("\n")
30
31 # Proporciones con totales.
32 proporciones <- prop.table(contingencia)
33 proporciones <- addmargins(proporciones)
34 cat("Tabla de contingencia con proporciones totales:\n")
35 print(proporciones)
36 cat("\n")

```

Aunque no ocurre con frecuencia, podríamos necesitar una matriz de confusión para más de dos variables. Veamos ahora un ejemplo con tres variables: Motor, Cambios y Transmisión. Para ello, tomamos una de las variables (en este caso, Motor) y creamos una subtabla por cada uno de sus niveles. Cada subtabla muestra las frecuencias para la combinación de las dos variables restantes cuando Motor tiene el nivel correspondiente, como muestra la tabla 2.7. En R, podemos obtener estas tablas como muestra el script 2.6. Desde luego, esta misma idea puede extenderse para cuatro o más variables categóricas.

Motor = Recto		Cambios		
		3 cambios	4 cambios	5 cambios
Transmision	Automático	3	4	0
	Manual	0	6	1

Motor = V		Cambios		
		3 cambios	4 cambios	5 cambios
Transmision	Automático	12	0	0
	Manual	0	2	4

Tabla 2.7: tabla de contingencia para tres variables.

Script 2.6: matriz de confusión para tres variables.

```

1 # Cargar datos.
2 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
3                   row.names = 1)
4

```

```

5 # Convertir la variable Cambios en categórica.
6 datos[["Cambios"]] <- factor(datos[["Cambios"]])
7
8 # Crear tabla de contingencia para las variables Transmision,
9 # Cambios y Motor.
10 contingencia <- ftable(datos[["Transmision"]], datos[["Cambios"]],
11                        datos[["Motor"]])
12
13 cat("Tabla de contingencia generada con ftable():\n")
14 print(contingencia)
15 cat("\n")
16
17 # Otra forma de crear la misma tabla.
18 xtabs(~ Cambios + Transmision + Motor, data = datos)
19 cat("Tabla de contingencia generada con xtabs():\n")
20 print(contingencia)
21 cat("\n")

```

### 2.1.3 Trabajando con datos agrupados

A menudo nos veremos en la necesidad de obtener estadísticas descriptivas de una variable separando las observaciones en grupos de acuerdo a una variable categórica. Para ello, el paquete `dplyr` ofrece la función `group_by()`, que podemos usar en conjunto con `summarise()`, como muestra el script 2.7. En dicho script, primero se agrupan las observaciones de acuerdo a la variable `Cambios`, y luego se efectúa una llamada a `summarise()` donde el primer argumento cuenta la cantidad de observaciones en el grupo actual y los argumentos restantes (que pueden ser tantos como se desee) corresponden a diferentes estadísticas descriptivas.

Script 2.7: estadísticas descriptivas para datos agrupados.

```

1 library(dplyr)
2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                  row.names = 1)
6
7 resumen <- group_by(datos, Cambios) %>%
8   summarise(count = n(), mean(Rendimiento), median(Rendimiento),
9             sd(Rendimiento), IQR(Rendimiento), mean(Potencia))
10
11 print(resumen)

```

## 2.2 REPRESENTACIÓN GRÁFICA DE DATOS

En esta sección revisaremos diversos tipos de gráficos que resultan útiles al momento de estudiar un conjunto de datos disponibles, considerando su definición, su utilidad y cómo se construyen en R. Para crear gráficos en R usaremos el paquete `ggpubr`. Algunos de los principales parámetros que usaremos para crear y editar gráficos con este paquete son:

- **data**: un data frame.
- **x**: string con el nombre de la variable  $x$ .
- **y**: string(s) con el(los) nombre(s) de la(s) variable(s) a graficar.
- **color**: color de delineado.
- **fill**: color de relleno.
- **palette**: paleta de colores cuando existen múltiples grupos.
- **linetype**: tipo de línea a emplear.
- **add**: permite agregar elementos adicionales al gráfico, como barras de error o la media, entre otros.
- **title**: título del gráfico.
- **xlab**: rótulo del eje  $x$ . Puede ocultarse usando `xlab = FALSE`.
- **ylab**: rótulo del eje  $y$ . Puede ocultarse usando `ylab = FALSE`.

### 2.2.1 Una variable numérica

El **histograma** resulta muy útil si queremos representar una única variable numérica y la muestra es grande. Podemos decir que este gráfico muestra una aproximación a la **densidad** (o distribución de frecuencias) para la variable, para lo que tenemos que dividir el rango de valores posibles en intervalos (generalmente iguales) y luego contar la cantidad de observaciones en cada intervalo. Para construir el gráfico, creamos una barra por cada intervalo, cuya altura (o longitud) es proporcional a la cantidad de observaciones en el intervalo representado. La figura 2.2 muestra histogramas creados con el script 2.8.

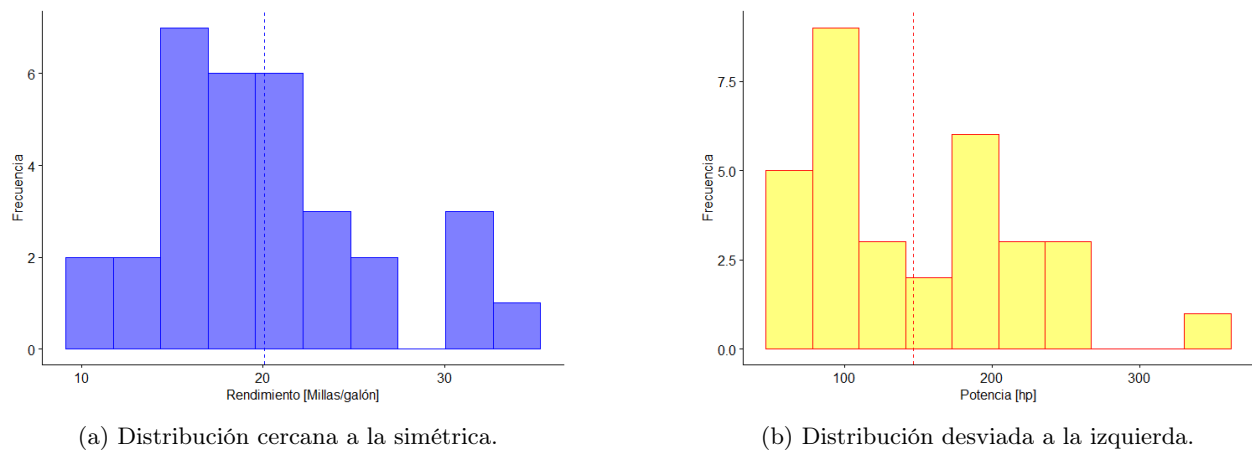


Figura 2.2: dos histogramas.

Script 2.8: histogramas para las variables Rendimiento y Potencia.

```
1 library(ggpubr)
2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                   row.names = 1)
6
7 # Histograma para la variable Rendimiento.
8 g1 <- gghistogram(datos,
9                  x = "Rendimiento",
10                 bins = 10,
11                 add = "mean",
```

```

12         xlab = "Rendimiento [Millas/galón]",
13         ylab = "Frecuencia",
14         color = "blue",
15         fill = "blue")
16
17 print(g1)
18
19 # Histograma para la variable Potencia.
20 g2 <- gghistogram(datos,
21                   x = "Potencia",
22                   bins = 10,
23                   add = "mean",
24                   xlab = "Potencia [hp]",
25                   ylab = "Frecuencia",
26                   color = "red",
27                   fill = "yellow")
28
29 print(g2)

```

A medida que avancemos en este libro, veremos que es muy importante conocer la **distribución de frecuencias** de una variable. Al observar la figura 2.2b, podemos ver que la frecuencia es mayor para potencias más bajas, pues las barras de la izquierda del gráfico son, en general, algo más altas que las de la derecha. Podría decirse que las observaciones se concentran a la izquierda y que hay una cola que se prolonga hacia la derecha. Cuando esto ocurre, decimos que la distribución está **desviada a la izquierda**, o que hay **asimetría negativa**. Análogamente, podría darse que la distribución estuviese desviada a la derecha o, equivalentemente, que presenta asimetría positiva. En el caso de la figura 2.2a, el histograma es más **simétrico**, pues las observaciones se aglomeran hacia el centro y hay colas tanto a la izquierda como a la derecha. Para ilustrar mejor la idea de la simetría, podemos revisar una vez más la figura 2.1, donde la población central es perfectamente simétrica y la inferior presenta asimetría positiva.

Otra ventaja de los histogramas es que permiten identificar modas de una variable, las cuales corresponden a barras que sean más prominentes que las de su entorno. Ambos ejemplos de la figura 2.2 son bimodales, pues tienen dos modas claramente identificables. Si bien es cierto que en ambos casos hay un único valor más frecuente (moda), podemos ver apreciar que existen dos “cumbres” o máximos locales.

Otro gráfico que usaremos a menudo es el de **gráfico de caja**. Es muy útil, pues su construcción considera 5 estadísticos para representar el conjunto de datos y además facilita la identificación de datos atípicos. La figura 2.3 muestra este gráfico para la variable *Potencia*, el cual fue creado con el script 2.9.

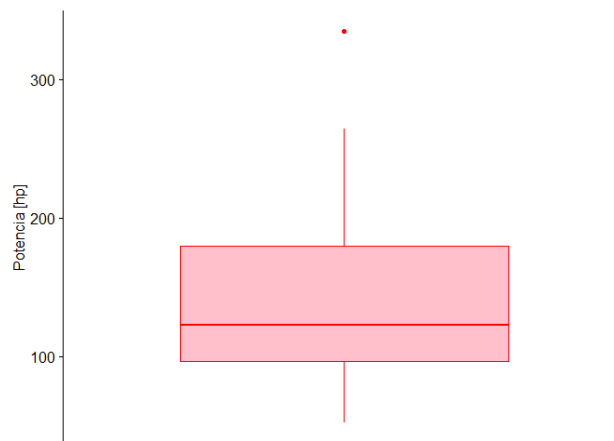


Figura 2.3: gráfico de caja para la variable *Potencia*.

Los extremos inferior y superior del rectángulo o caja de la figura 2.3 corresponden, respectivamente, al

primer y al tercer cuartil, mientras que la línea horizontal al interior de la caja denota la mediana. Así, la caja engloba el 50% central de los datos, y su altura corresponde al rango intercuartil. Las barras que se extienden por sobre y por debajo de la caja, llamadas bigotes, capturan aquellos datos fuera de la caja central y que estén situados a no más de 1,5 veces el IQR. Cualquier observación que esté más allá de la caja y los bigotes se representa como un punto, el cual podría tratarse de una observación atípica.

Script 2.9: gráfico de caja para la variable **Potencia**.

```
1 library(ggpubr)
2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                   row.names = 1)
6
7 g <- ggboxplot(datos[["Potencia"]],
8               color = "red",
9               fill = "pink",
10              ylab = "Potencia [hp]")
11
12 g <- g + rremove("x.ticks")
13 g <- g + rremove("x.text")
14 g <- g + rremove("x.title")
15
16 print(g)
```

### 2.2.2 Una variable categórica

Si queremos representar una única variable categórica, lo más adecuado es usar un **gráfico de barras**, pues cada barra es tan larga como la proporción de valores presentes en cada nivel de la variable. La figura 2.4 muestra el gráfico de barras correspondiente a la tabla 2.2, elaborado mediante el script 2.10.

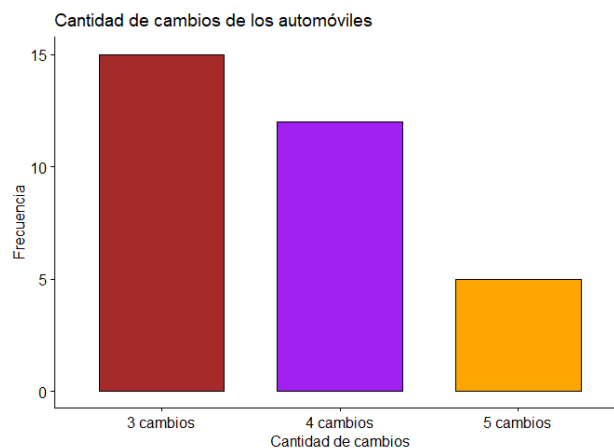


Figura 2.4: gráfico de barras para la variable **Cambios**.

Script 2.10: gráfico de barras para la variable **Cambios**.

```
1 library(ggpubr)
2
```

```

3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                   row.names = 1)
6
7 # Crear la tabla de frecuencias para la variable Cambios y convertirla a
8 # data frame.
9 contingencia <- as.data.frame(xtabs(~ Cambios, data = datos))
10
11 # Crear el gráfico de barras.
12 g <- ggbarplot(contingencia,
13               x = "Cambios",
14               y = "Freq",
15               fill = c("brown", "purple", "orange"),
16               title = "Cantidad de cambios de los automóviles",
17               xlab = "Cantidad de cambios",
18               ylab = "Frecuencia")
19
20 print(g)

```

Otra alternativa para representar una única variable categórica es el **gráfico de torta**, que se presenta en la figura 2.5 y se construye en R como muestra el script 2.11.

Cantidad de cambios de los automóviles

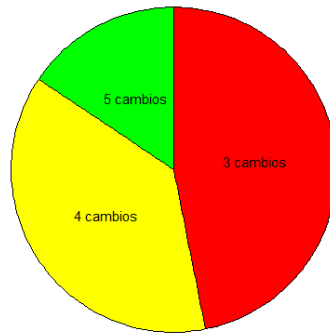


Figura 2.5: gráfico de torta para la variable Cambios.

Script 2.11: gráfico de torta para la variable Cambios.

```

1 library(ggpubr)
2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                   row.names = 1)
6
7 # Crear la tabla de frecuencias y convertirla a data frame.
8 contingencia <- as.data.frame(xtabs(~ Cambios, data = datos))
9
10 # Crear gráfico de torta.
11 g <- ggpie(contingencia,
12           x = "Freq",
13           label = "Cambios",
14           fill = c("red", "yellow", "green"),
15           title = "Cantidad de cambios de los automóviles",
16           lab.pos = "in")
17

```

```
18 print(g)
```

### 2.2.3 Dos variables numéricas

Los **gráficos de dispersión** son adecuados en este caso. Se caracterizan porque muestran información caso a caso, ya que cada punto del gráfico corresponde a una observación. Por ejemplo, el gráfico de la figura 2.6, creado mediante el script 2.12, muestra este tipo de gráfico para las variables **Rendimiento** y **Peso**.

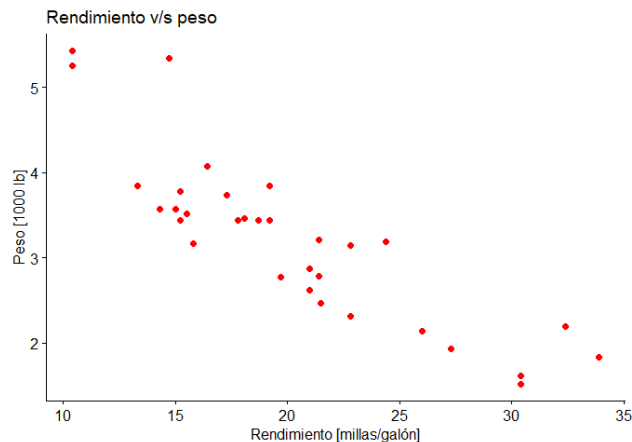


Figura 2.6: gráfico de dispersión para las variables **Rendimiento** y **Peso**.

Script 2.12: gráfico de dispersión para las variables **Rendimiento** y **Peso**.

```
1 library(ggpubr)
2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                   row.names = 1)
6
7 # Crear gráfico de dispersión.
8 g <- ggscatter(datos,
9               x = "Rendimiento",
10              y = "Peso",
11              color = "red",
12              title = "Rendimiento v/s peso",
13              xlab = "Rendimiento [millas/galón]",
14              ylab = "Peso [1000 lb]")
15
16 print(g)
```

Los gráficos de dispersión también son muy útiles para identificar si dos (o más) variables están relacionadas. La figura 2.7 (creada mediante el script 2.13) muestra tres gráficos de dispersión diferentes: en el de la izquierda, se aprecia que las variables **Peso** y **Cuarto\_milla** son independientes, pues no hay una tendencia definida en la organización de los puntos. En el gráfico del centro, en cambio, podemos ver que la potencia tiende a aumentar a medida que también lo hace el peso, por lo que ambas variables están positivamente asociadas. Por último, el gráfico de la derecha nos muestra que las variables **Peso** y **Rendimiento** presentan asociación negativa, puesto que a medida que la primera aumenta, la segunda disminuye.



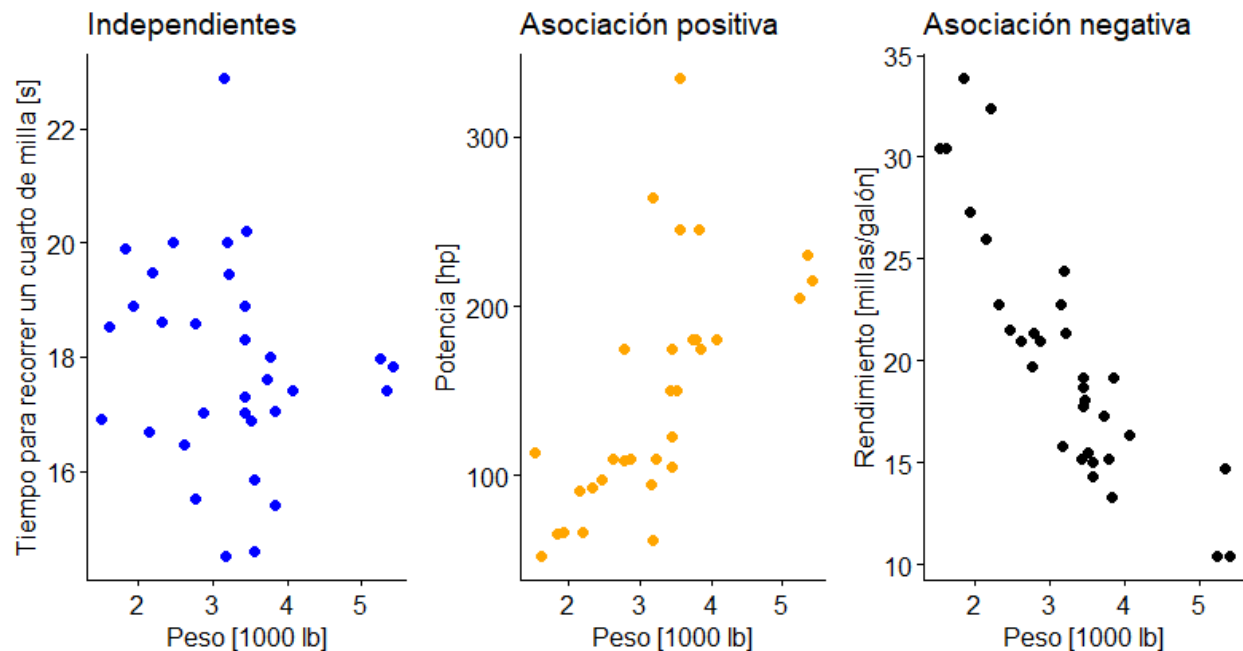


Figura 2.7: gráficos de dispersión con diferentes tipos de asociación entre las variables.

Script 2.13: gráficos de dispersión con diferentes tipos de asociación entre las variables.

```

1 library(ggpubr)
2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                   row.names = 1)
6
7 # Gráfico para variables independientes.
8 g1 <- ggscatter(datos,
9                 x = "Peso",
10                 y = "Cuarto_milla",
11                 color = "blue",
12                 title = "Independientes",
13                 xlab = "Peso [1000 lb]",
14                 ylab = "Tiempo para recorrer un cuarto de milla [s]")
15
16 # Gráfico para variables con asociación positiva.
17 g2 <- ggscatter(datos,
18                 x = "Peso",
19                 y = "Potencia",
20                 color = "orange",
21                 title = "Asociación positiva",
22                 xlab = "Peso [1000 lb]",
23                 ylab = "Potencia [hp]")
24
25 # Gráfico para variables con asociación negativa.
26 g3 <- ggscatter(datos,
27                 x = "Peso",
28                 y = "Rendimiento",
29                 color = "black",
30                 title = "Asociación negativa",
31                 xlab = "Peso [1000 lb]",

```

```

32         ylab = "Rendimiento [millas/galón]")
33
34 # Crear figura con tres gráficos.
35 g <- ggarrange(g1 ,g2 ,g3, ncol = 3, nrow = 1, common.legend = TRUE)
36
37 print(g)

```

## 2.2.4 Dos variables categóricas

Similares al gráfico de barras para una variable categórica, los **gráficos de barras apiladas, agrupadas y estandarizadas** permiten visualizar la matriz de confusión entre dos variables y encontrar posibles relaciones entre ellas. La figura 2.8, creada con el script 2.14, ejemplifica esta familia de gráficos usando para ello las variables **Cambios** y **Motor**.

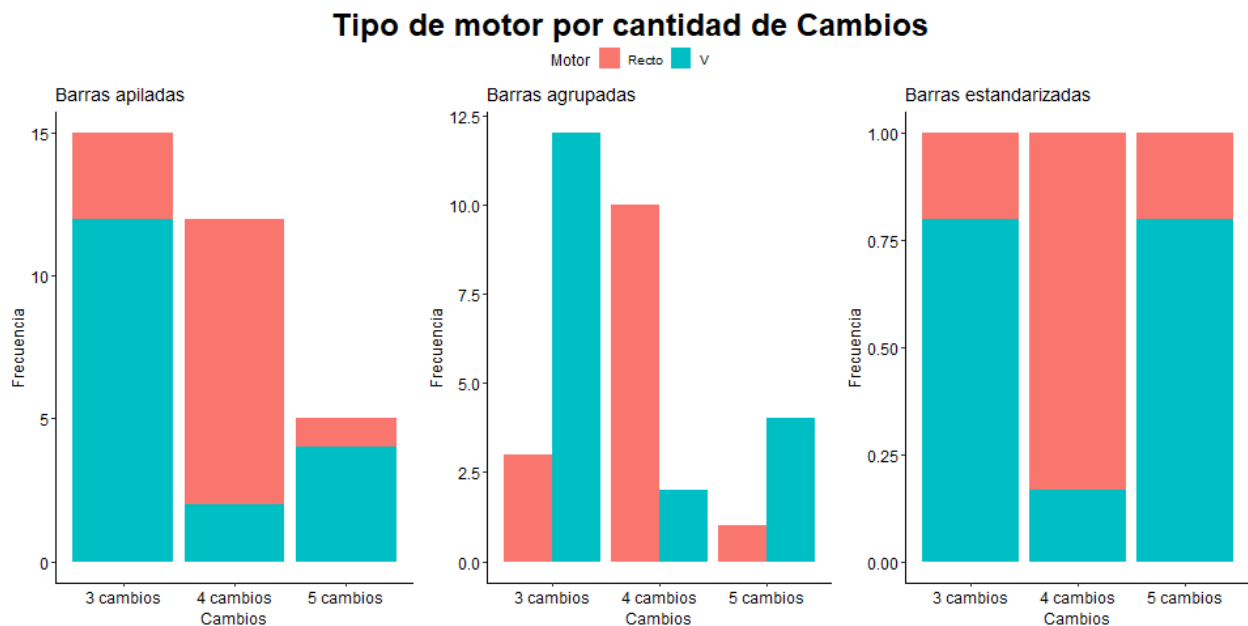


Figura 2.8: gráficos de barras para las variables **Cambios** y **Motor**.

El **gráfico de barras apiladas**, a la izquierda en la figura 2.8, muestra tres barras cuya altura corresponde a la frecuencia de la cantidad de cambios, al igual que en la figura 2.4, pero ahora cada barra está subdividida en secciones de distinto color para cada tipo de motor. La altura de cada sección está dada por la frecuencia del tipo de motor para la cantidad de cambios representada en la barra.

Similar al anterior, el gráfico de la derecha en la figura 2.8, que corresponde a un **gráfico de barras estandarizadas**, muestra barras de igual altura para cada cantidad de cambios representando claramente los cambios en la proporción de cada tipo de motor por la cantidad de cambios. Se puede apreciar que los automóviles con 3 y 5 cambios tienen mayoritariamente motores en forma de V, ambas en igual proporción, mientras que el uso de motores rectos se da principalmente en automóviles de 4 cambios.

El **gráfico de barras agrupadas**, al centro en la figura 2.8, es equivalente al de la izquierda, pero en lugar de dividir una barra en segmentos, muestra barras contiguas para cada tipo de motor.

Script 2.14: gráficos de barras para las variables **Cambios** y **Motor**.

```

1 library(ggpubr)
2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                   row.names = 1)
6
7 # Crear tabla de contingencia para las variables Motor y Cambios,
8 # y guardarla como data frame.
9 tabla <- xtabs(~ Motor + Cambios, data = datos)
10 contingencia <- as.data.frame(tabla)
11
12 # Crear gráfico de barras segmentadas.
13 g1 <- ggplot(contingencia, aes(fill = Motor, y = Freq, x = Cambios))
14 g1 <- g1 + geom_bar(position = "stack", stat = "identity")
15 g1 <- g1 + labs(y = "Frecuencia") + ggtitle("Barras apiladas")
16 g1 <- g1 + theme_pubr()
17
18 # Crear gráfico de barras agrupadas.
19 g2 <- ggplot(contingencia, aes(fill = Motor, y = Freq, x = Cambios))
20 g2 <- g2 + geom_bar(position = "dodge", stat = "identity")
21 g2 <- g2 + labs(y = "Frecuencia") + ggtitle("Barras agrupadas")
22 g2 <- g2 + theme_pubr()
23
24 # # Crear gráfico de barras segmentadas estandarizado.
25 g3 <- ggplot(contingencia, aes(fill = Motor, y = Freq, x = Cambios))
26 g3 <- g3 + geom_bar(position = "fill", stat = "identity")
27 g3 <- g3 + labs(y = "Frecuencia") + ggtitle("Barras estandarizadas")
28 g3 <- g3 + theme_pubr()
29
30 # Crear una figura que contenga los tres gráficos.
31 g <- ggarrange(g1, g2, g3, nrow = 1, common.legend = TRUE)
32
33 # Agregar un título común en negrita y con fuente de 24 puntos.
34 titulo <- text_grob("Tipo de motor por cantidad de Cambios",
35                    face = "bold", size = 24)
36
37 g <- annotate_figure(g, top = titulo)
38
39 # Guardar la figura en formato png con tamaño 960 x 480 pixeles.
40 ggexport(g, filename = "C:/Inferencia/f-barras-2.png",
41          height = 480, width = 960)

```

Similar al gráfico de barras para dos variables, el **gráfico de mosaico** permite representar una tabla de contingencia. Para ello, divide un área en regiones y el área de cada región es proporcional al porcentaje de observaciones que representa. La figura 2.9, creada con el script 2.15 ejemplifica el uso de este tipo de gráficos, usando para ello las variables Cambios y Motor. En ella, el ancho de cada columna es proporcional a la cantidad de automóviles que tienen la correspondiente cantidad de cambios, mientras que la altura de cada barra de las columnas refleja la proporción de automóviles con un determinado tipo de motor.

Si nos fijamos bien en la figura 2.9, podemos ver claramente que los vehículos con 5 cambios son, por mucho, los menos frecuentes y que los con 3 cambios son algo más frecuentes que los que tienen 4 cambios. Del mismo modo, podemos ver que, para vehículos de 3 y 5 cambios, la proporción de vehículos con motor recto es la misma, y mucho menor que la de aquellos con motor en forma de V. Sin embargo, este último no es muy frecuente en automóviles con 4 cambios.

Cabe destacar que, para este tipo de gráfico, se requiere emplear el paquete **ggmosaic**.

Script 2.15: gráfico de mosaico para las variables Cambios y Motor.

```

1 library(ggmosaic)

```

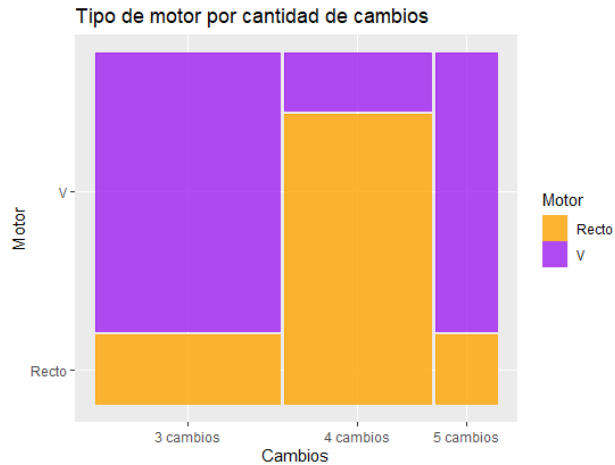


Figura 2.9: gráfico de mosaico para las variables Cambios y Motor.

```

2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                   row.names = 1)
6
7 # Crear tabla de contingencia para las variables gear y vs,
8 # y guardarla como data frame.
9 tabla <- xtabs(~ Cambios + Motor, data = datos)
10 contingencia <- as.data.frame(tabla)
11
12 # Crear gráfico de mosaico.
13 g <- ggplot(data = contingencia)
14 g <- g + geom_mosaic(aes(weight = Freq, x = product(Cambios), fill = Motor))
15
16 g <- g + labs(y = "Motor", x = "Cambios",
17              title = "Tipo de motor por cantidad de cambios")
18
19 g <- g + scale_fill_manual(values=c("orange", "purple"))
20
21 print(g)

```

### 2.2.5 Una variable numérica y otra categórica

Desde luego, también es importante poder comparar diferentes grupos de observaciones de acuerdo a una característica categórica, para lo cual los gráficos pueden ser de gran ayuda. Por ejemplo, la figura 2.10, creada mediante el script 2.16 muestra un gráfico de cajas para la variable Rendimiento agrupada por el número de cambios de los automóviles.

Script 2.16: gráfico de cajas por grupo.

```

1 library(ggpubr)
2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,

```

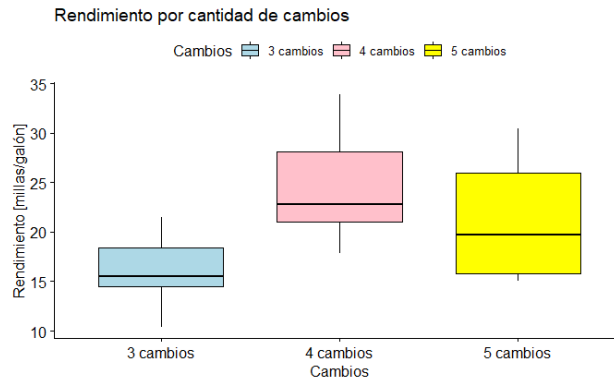


Figura 2.10: gráfico de cajas por grupo.

```

5         row.names = 1)
6
7 g <- ggboxplot(datos, x = "Cambios",
8               y = "Rendimiento",
9               palette = c("light blue", "pink", "yellow"),
10              fill = "Cambios",
11              title = "Rendimiento por cantidad de cambios",
12              xlab = "Cambios",
13              ylab = "Rendimiento [millas/galón]")
14
15 print(g)

```

Una buena alternativa, si la cantidad de observaciones es pequeña, es el **gráfico de tiras**, similar al gráfico de dispersión. El script 2.17 construye este gráfico para la variable **Rendimiento** agrupada según los niveles de la variable **Cambios**, obteniéndose como resultado la figura 2.11.

Script 2.17: gráfico de tiras.

```

1 library(ggpubr)
2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                  row.names = 1)
6
7 g <- ggstripchart(datos, x = "Cambios",
8                  y = "Rendimiento",
9                  palette = c("blue", "red", "dark green"),
10                 color = "Cambios",
11                 title = "Rendimiento por cantidad de cambios",
12                 xlab = "Cambios",
13                 ylab = "Rendimiento [millas/galón]")
14
15 print(g)

```

## 2.3 EJERCICIOS PROPUESTOS

1. ¿Cuándo es apropiado utilizar un gráfico de puntos para revisar datos?

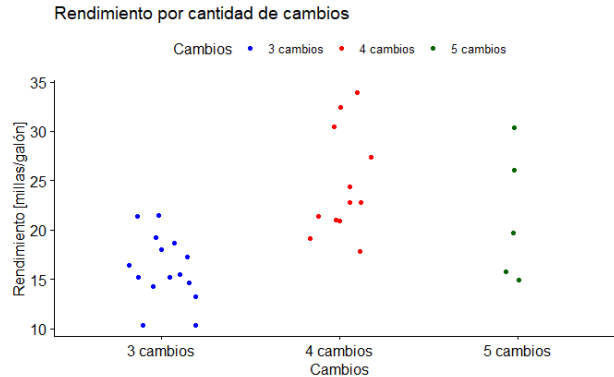
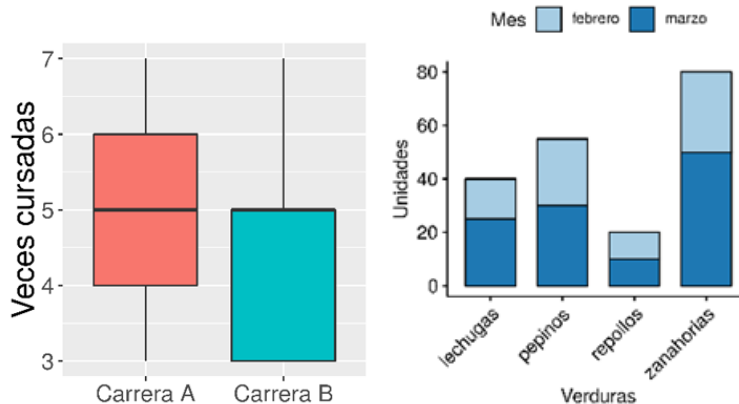


Figura 2.11: gráfico de tiras.

2. ¿Cuándo la mediana caracteriza mejor a un conjunto de datos que la media?
3. Da ejemplos de tres variables que posiblemente tengan una distribución simétrica, con asimetría positiva y con asimetría negativa, respectivamente. Justifique bien cada caso.
4. Describe un estudio en que posiblemente los datos recolectados tengan una distribución bimodal.
5. ¿Qué tipo de información buscada llevaría a utilizar un gráfico de dispersión?
6. ¿Por qué es importante conocer una medida de dispersión (variabilidad) de un conjunto de datos? Dé un ejemplo para clarificar su respuesta.
7. Considera la representación de la figura 2.12a de los datos obtenidos al tomar muestras aleatorias de estudiantes de dos carreras de la Facultad de Ingeniería para estudiar si el número de veces que se cursan las tres asignaturas de física en el Módulo Básico de Ingeniería depende de la carrera de los alumnos. Compara las distribuciones de ambos grupos. ¿En qué se parecen y en qué se diferencian?



(a) Ejercicio 7.

(b) Ejercicio 8.

Figura 2.12: gráficos para los ejercicios propuestos.

8. El gráfico de la figura 2.12b muestra las unidades de verduras vendidas en uno de los kioscos de la Universidad durante los meses anteriores. Construye la tabla de contingencia correspondiente a los datos que se representan. ¿Qué mes tuvo mayores ventas? ¿En qué proporción?
9. ¿Cómo se puede generar la secuencia 0.00, 0.25, 0.50, ..., 2.75, 3.00 en R?
10. Resuelve en R los siguientes ejercicios. Considera para ello el conjunto de datos nativo de R `chickwts`.
  - a) ¿Qué son los cuartiles y cómo se pueden obtener para los pesos de los pollitos reportados en la columna `weight`?
  - b) ¿Cómo obtener los cuartiles del ejercicio anterior por cada tipo de alimento en la columna `feed`?
  - c) ¿Cómo obtener un histograma de los pesos de los pollitos?

- d)* ¿Cómo se obtiene un gráfico de cajas para comparar los pesos de los pollitos por tipo de alimento suministrado?
11. Investiga acerca del uso del paquete de R `ggplot2` para la creación de gráficos.





## CAPÍTULO 3. VARIABLES ALEATORIAS Y DISTRIBUCIONES DE PROBABILIDAD

Los conceptos que estudiaremos en este capítulo pueden resultar algo difíciles de entender, por lo que si necesitas más material, puedes consultar las fuentes en que se basa este capítulo: Diez y col. (2017, pp. 104-157) y Freund y Wilson (2003, pp. 104-106).

Definimos como **variable aleatoria** una variable o un proceso cuyo resultado sea numérico. Dichas variables se nombran con letras mayúsculas, y denotamos sus posibles valores por la letra minúscula correspondiente, acompañada de un subíndice. Las variables aleatorias tienen una **distribución de probabilidad**, la cual define la probabilidad de que ocurran los diferentes valores que dicha variable puede tomar.

### 3.1 VARIABLES ALEATORIAS

La definición de **variable aleatoria continua** es en realidad bastante sencilla: es una variable que puede tomar cualquiera de los infinitos valores posibles dentro de un intervalo.

Una **variable aleatoria discreta**, en cambio, solo puede tomar un conjunto finito de valores. Un ejemplo típico de variable aleatoria puede ser el lanzamiento de un dado. Si el dado está bien balanceado, tendremos igual probabilidad de obtener cualquiera de las caras. Pero es sabido que algunos tramposos fabrican dados adulterados para favorecer, por ejemplo, la obtención de valores 1 y 6. Una distribución aleatoria de la variable lanzamiento de un dado adulterado ( $X$ ) podría ser la que se presenta en la tabla 3.1.

$i$	1	2	3	4	5	6	Total
$x_i$	1	2	3	4	5	6	-
$P(X = x_i)$	0.250	0.125	0.125	0.125	0.125	0.250	1.000

Tabla 3.1: distribución de probabilidad para el lanzamiento de un dado adulterado.

El **valor esperado**, denotado como  $E(X)$  o  $\mu$ , corresponde al resultado promedio de una variable aleatoria. Para una variable aleatoria discreta, se calcula sumando los valores posibles ponderados por su probabilidad, como muestra la ecuación 3.1.

$$E(X) \equiv \mu = \sum_{i=1}^n x_i P(X = x_i) \quad (3.1)$$

También podemos calcular qué tan alejado podría estar el valor obtenido del valor esperado por medio de la varianza general, denotada por  $Var(X)$  o  $\sigma^2$ , que se calcula como la media de los cuadrados de la diferencia con respecto a la media ponderada según la probabilidad de ocurrencia, como muestra la ecuación 3.2. Una vez más, la desviación estándar corresponde a la raíz cuadrada de la varianza.

$$Var(X) \equiv \sigma^2 = \sum_{i=1}^n (x_i - \mu)^2 P(X = x_i) \quad (3.2)$$

En R, el paquete `DiscreteRV` permite trabajar con variables aleatorias discretas, como se ejemplifica en el script 3.1 (Cross, 2017).

Script 3.1: variables aleatorias discretas en R.

```
1 library(discreteRV)
2
3 # Crear una variable discreta para representar el dado
4 # adulterado de la tabla 3.1.
5 resultados <- 1:6
6 probabilidades = c(0.25, 0.125, 0.125, 0.125, 0.125, 0.25)
7 X <- RV(outcomes = resultados, probs = probabilidades)
8
9 # Calcular el valor esperado.
10 esperado <- E(X)
11 cat("Valor esperado:", esperado, "\n")
12
13 # Calcular la varianza.
14 varianza <- V(X)
15 cat("Varianza:", varianza, "\n")
16
17 # Calcular la desviación estándar.
18 desviacion <- SD(X)
19 cat("Desviación estándar:", desviacion, "\n")
```

Para ayudarnos a entender mejor la noción de distribución de probabilidad, veamos la figura 3.1 (obtenida mediante el script 3.2). Ella nos muestra, de izquierda a derecha, las distribuciones de probabilidad para el puntaje total obtenido al lanzar 5, 10 y 20 dados, respectivamente.

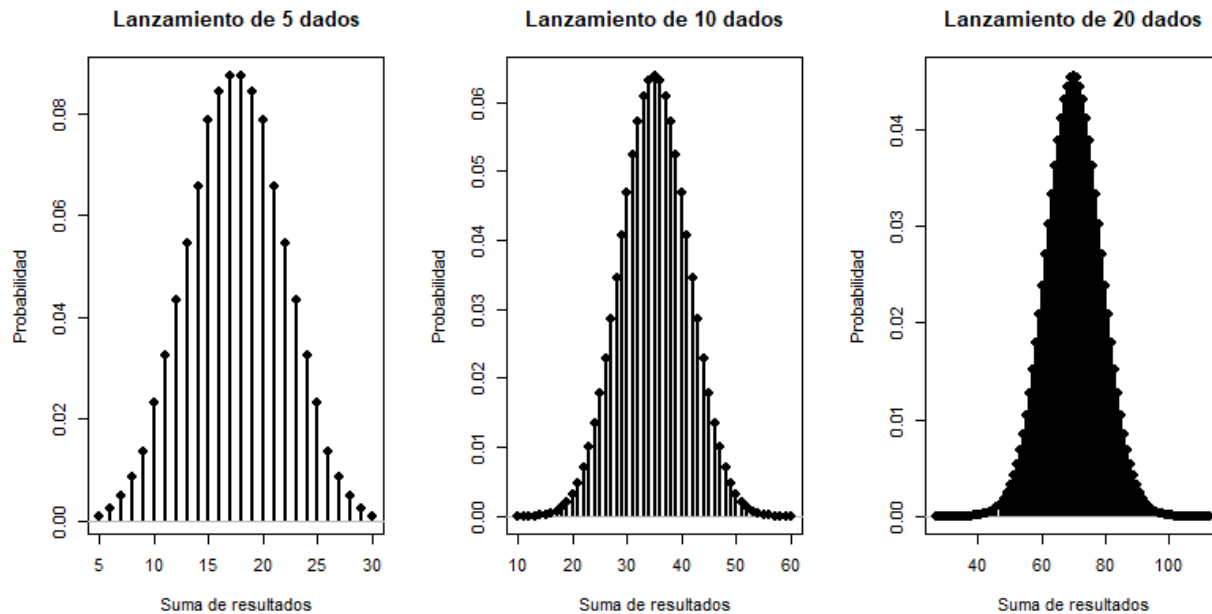


Figura 3.1: distribución de probabilidad para varios lanzamientos de un dado cargado.

Script 3.2: histogramas de variables aleatorias discretas en R.

```
1 library(discreteRV)
2 library(ggpubr)
3
4 # Crear una variable discreta para representar el dado
5 # adulterado de la tabla 4.1.
6 resultados <- 1:6
```

```

7 probabilidades = c(0.25, 0.125, 0.125, 0.125, 0.125, 0.25)
8 X <- RV(outcomes = resultados, probs = probabilidades)
9
10 # Crear vector con los resultados de 5 lanzamientos del dado.
11 lanzar_5 <- SofIID(X, n=5)
12
13 # Crear vector con los resultados de 10 lanzamientos del dado.
14 lanzar_10 <- SofIID(X, n=10)
15
16 # Crear vector con los resultados de 20 lanzamientos del dado.
17 lanzar_20 <- SofIID(X, n=20)
18
19 # Graficar los resultados.
20 par(mfrow=c(1, 3))
21
22 plot(lanzar_5,
23      main="Lanzamiento de 5 dados",
24      xlab="Suma de resultados",
25      ylab="Probabilidad")
26
27 plot(lanzar_10,
28      main="Lanzamiento de 10 dados",
29      xlab="Suma de resultados",
30      ylab="Probabilidad")
31
32 plot(lanzar_20,
33      main="Lanzamiento de 20 dados",
34      xlab="Suma de resultados",
35      ylab="Probabilidad")

```

Conocer la distribución de probabilidad de una variable discreta nos ayuda a hacer estimaciones útiles. A modo de ejemplo, supongamos que un ingeniero de software debe crear un programa que resuelva un problema (siempre con instancias del mismo tamaño) con un tiempo de respuesta no mayor a 25 segundos. El histograma de la figura 3.2 muestra los tiempos de ejecución obtenidos para 500 pruebas de la solución propuesta, donde se observa que 30 de ellas tardaron en realidad más de 25 segundos, con un rango que va de 0 a 30 segundos. Así, podemos estimar la probabilidad de que el tiempo de ejecución sea mayor a 25 segundos dividiendo la cantidad de observaciones que cumplen este criterio por la cantidad total de instancias, como muestra la ecuación 3.3.

$$P(X > 25) = \frac{30}{500} = 0.06 \quad (3.3)$$

Frecuentemente resulta más adecuado expresar o modelar un fenómeno como una combinación de dos o más variables aleatorias. Por ejemplo, un jugador de baloncesto puede anotar canastas de uno, dos o tres puntos dependiendo de si encesta con un tiro libre, un lanzamiento desde dentro del área o desde fuera del área. Así, se tienen tres variables aleatorias:

1.  $X$ : Anotaciones por tiro libre.
2.  $Y$ : Anotaciones desde dentro del área.
3.  $Z$ : Anotaciones desde fuera del área.

Podemos representar el total de puntos anotados por el jugador como la suma de los puntos anotados de las tres formas posibles, lo que corresponde a una **combinación lineal** de las variables  $X$ ,  $Y$  y  $Z$ . La fórmula general de una combinación lineal de  $n$  variables está dada por la ecuación 3.4, donde cada  $X_i$  corresponde a una variable aleatoria y cada  $c_i$  es una constante conocida.

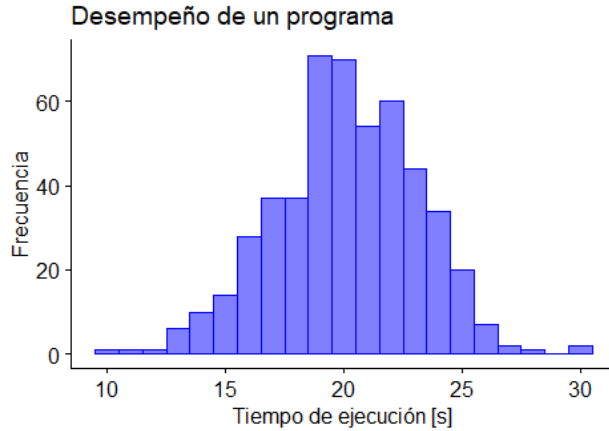


Figura 3.2: histograma para el desempeño del programa.

$$\sum_{i=1}^n c_i X_i \quad (3.4)$$

Cuando las variables de una combinación lineal son independientes<sup>1</sup>, podemos calcular el valor esperado y la varianza de la combinación lineal usando las ecuaciones 3.5 y 3.6. Una vez más, la desviación estándar está dada por la raíz cuadrada de la varianza.

$$E\left(\sum_{i=1}^n c_i X_i\right) = \sum_{i=1}^n c_i E(X_i) \quad (3.5)$$

$$Var\left(\sum_{i=1}^n c_i X_i\right) = \sum_{i=1}^n c_i^2 Var(X_i) \quad (3.6)$$

Por supuesto, en R también podemos trabajar con combinaciones lineales de variables aleatorias discretas, como muestra el script 3.3.

Script 3.3: combinación lineal de variables aleatorias discretas en R.

```

1 library(discreteRV)
2
3 # Crear una variable discreta para representar el dado adulterado de la tabla
4 # 3.1, y calcular su valor esperado, varianza y desviación estándar.
5 resultados <- 1:6
6 probabilidades = c(0.25, 0.125, 0.125, 0.125, 0.125, 0.25)
7 X <- RV(outcomes = resultados, probs = probabilidades)
8 esperado_x <- E(X)
9 varianza_x <- V(X)
10 desviacion_x <- SD(X)
11 cat("E(X):", esperado_x, "\n")
12 cat("V(X):", varianza_x, "\n")
13 cat("SD(X):", desviacion_x, "\n\n")
14
15 # Crear una variable aleatoria para un dado balanceado, y calcular su valor
16 # esperado, varianza y desviación estándar.
17 Y <- RV(outcomes = resultados, probs = 1/6)
```

<sup>1</sup>Si las variables no son independientes, se requieren métodos más complejos fuera del alcance de este libro.

```

18 esperado_y <- E(Y)
19 varianza_y <- V(Y)
20 desviacion_y <- SD(Y)
21 cat("E(Y):", esperado_y, "\n")
22 cat("V(Y):", varianza_y, "\n")
23 cat("SD(Y):", desviacion_y, "\n\n")
24
25 # Crear una combinación lineal de variables aleatorias, y calcular su valor
26 # esperado, varianza y desviación estándar.
27 Z <- 0.5 * X + 0.5 * Y
28 esperado_z <- E(Z)
29 varianza_z <- V(Z)
30 desviacion_z <- SD(Z)
31 cat("E(Z):", esperado_z, "\n")
32 cat("V(Z):", varianza_z, "\n")
33 cat("SD(Z):", desviacion_z)

```

Al examinar con mayor detención los gráficos de la figura 3.1 podemos apreciar que, a medida que se efectúan más lanzamientos del dado, el histograma se asemeja cada vez más a una curva continua, la cual recibe el nombre de **función de densidad de probabilidad**, o simplemente **distribución** o **densidad**.

Las distribuciones tienen la propiedad de que el área total bajo la curva siempre es 1, lo que resulta muy útil al momento de calcular probabilidades, pues basta con calcular el área bajo la curva del segmento deseado. Volviendo al ejemplo del desempeño del programa, presentado en la figura 3.2, el tiempo de ejecución es en realidad una variable continua. Así, la probabilidad de que el tiempo de ejecución sea mayor a 25 segundos corresponde al área coloreada en el gráfico de la figura 3.3, con un valor de 0,048<sup>2</sup>.

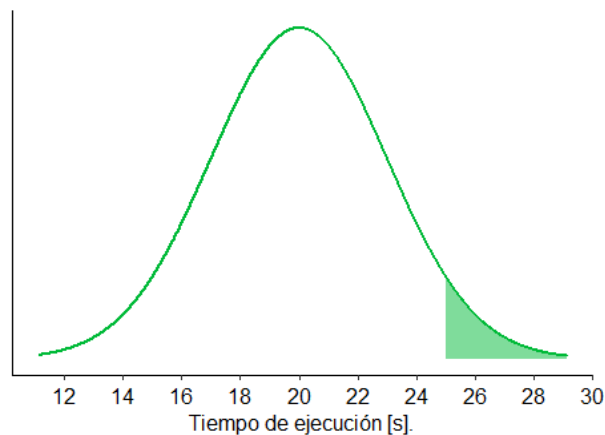


Figura 3.3: distribución para el desempeño del programa.

## 3.2 DISTRIBUCIONES CONTINUAS

Existen múltiples funciones de distribución continua que son de uso frecuente en estadística, las cuales se describen a continuación.

---

<sup>2</sup>El cálculo de esta probabilidad se aborda en el siguiente apartado

### 3.2.1 Distribución normal

También conocida como **distribución gaussiana**, la **distribución normal** es la más ampliamente empleada en estadística, pues muchas variables se acercan a esta distribución. Se caracteriza por ser unimodal y simétrica, con forma de campana. La figura 3.3 ejemplifica esta distribución.

La distribución normal se usa para modelar diversos fenómenos y podemos ajustarla mediante dos parámetros:

- $\mu$ : la media, que desplaza el centro de la curva a lo largo del eje x.
- $\sigma$ : la desviación estándar, que modifica qué tan dispersos están los datos con respecto a la media.

Así, denotamos este tipo de distribución por  $N(\mu, \sigma)$ . La figura 3.4, creada mediante el script 3.4, muestra dos ejemplos superpuestos de distribución normal:  $N(\mu = 0, \sigma = 1)$  en azul y  $N(\mu = 10, \sigma = 6)$  en rojo.

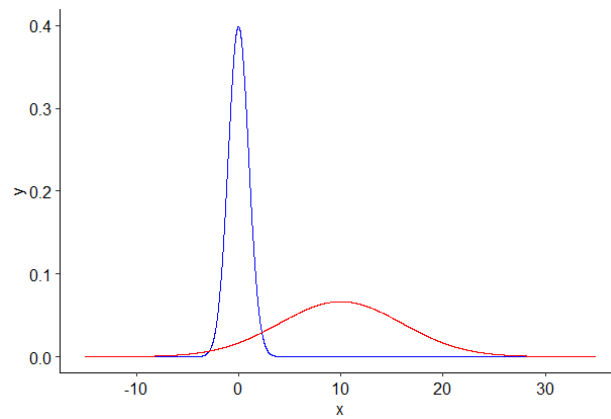


Figura 3.4: dos ejemplos superpuestos de distribución normal.

Script 3.4: graficando dos ejemplos de distribución normal.

```
1 library(ggpubr)
2
3 # Generar valores para una distribución normal con media 0 y
4 # desviación estándar 1.
5 media <- 0
6 desv_est <- 1
7 x <- seq(-15, 35, 0.01)
8 y <- dnorm(x, mean = media, sd = desv_est)
9 normal_1 <- data.frame(x, y)
10
11 # Repetir el proceso para una distribución normal con media 10
12 # y desviación estándar 6.
13 media <- 10
14 desv_est <- 6
15 x <- seq(-15, 35, 0.01)
16 y <- dnorm(x, mean = media, sd = desv_est)
17 normal_2 <- data.frame(x, y)
18
19 # Graficar ambas distribuciones.
20 g <- ggplot(normal_1, aes(x, y)) + geom_line(color = "blue")
21 g <- g + geom_line(data = normal_2, color = "red")
22 g <- g + theme_pubr()
23
24 print(g)
```

Antes de continuar, fijémonos en las líneas 8 y 16 del script 3.4, donde se usa la función `dnorm(x, mean, sd)`. Esta función calcula la densidad de una distribución normal. Además de `dnorm()`, R nos ofrece otras funciones que también resultan de mucha ayuda:

- `pnorm(q, mean, sd, lower.tail)`: permite encontrar percentiles, los cuales corresponden a la **función de distribución acumulada** (es decir, la probabilidad de que la variable tome valores menores o iguales que un valor dado), a partir de las probabilidades.
- `qnorm(p, mean, sd, lower.tail)`: encuentra el percentil para las probabilidades dadas en `p`, por lo que es la función inversa de `pnorm()`.
- `rnorm(n, mean, sd)`: genera aleatoriamente `n` observaciones de la distribución normal especificada.

Los argumentos de esta familia de funciones son:

- `x`, `q`: vector de cuantiles (percentiles).
- `p`: vector de probabilidades.
- `mean`: media de la distribución normal.
- `sd`: desviación estándar de la distribución normal.
- `lower.tail`: valor lógico que señala cuál de los dos extremos o colas de la distribución emplear.
- `n`: tamaño del vector resultante.

Es importante señalar que, por defecto, `lower.tail` toma el valor verdadero, con lo que `pnorm()` y `qnorm()` operan con la cola inferior de la distribución. Si, en cambio, `lower.tail = FALSE`, dichas funciones operan con la cola superior (es decir, `pnorm()` nos entrega la probabilidad de que la variable tome valores mayores que un valor dado).

Una **regla empírica** muy útil al momento de trabajar con distribuciones normales es la llamada regla 68-95-99.7, ilustrada en la figura 3.5, la cual establece que:

- Cerca de 68 % de las observaciones se encuentran a una distancia de una desviación estándar de la media.
- Alrededor de 95 % de las observaciones se encuentran a una distancia de dos desviación estándar de la media.
- Aproximadamente 99.7 % de las observaciones se encuentran a una distancia de tres desviación estándar de la media.

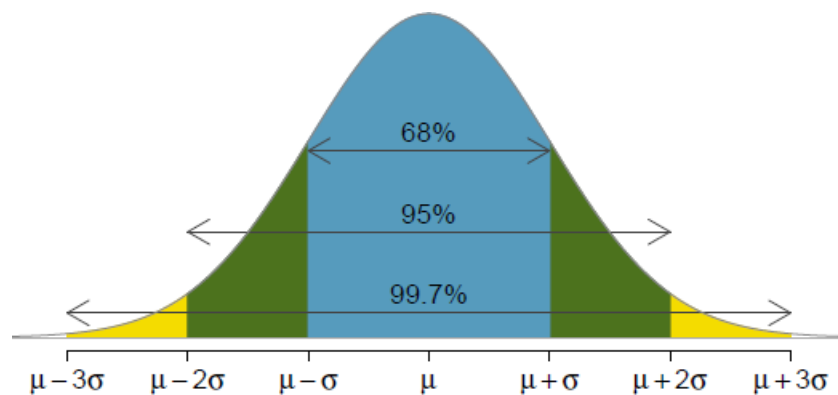


Figura 3.5: regla empírica de la distribución normal. Fuente: Diez y col. (2017, p. 136).

Muchas pruebas estadísticas operan bajo el supuesto de que los datos siguen una distribución normal. Como se insinuó en párrafos anteriores, la normalidad es siempre una aproximación, por lo que debemos verificar que el supuesto de una distribución normal sea aceptable. Una buena herramienta para ello es el **gráfico cuantil-cuantil**, también llamado **gráfico Q-Q**, que se muestra en la figura 3.6 y que podemos construir en R como muestra el script 3.5. En él podemos distinguir los siguientes elementos: un grupo de puntos, una recta y una región coloreada. Los puntos corresponden a las observaciones, mientras que la recta representa la distribución normal. En consecuencia, mientras más se asemeje el patrón que forman los puntos a la recta,

más parecida será la distribución a la normal. La banda coloreada establece el margen aceptable para suponer normalidad en el conjunto de datos. Así, para el conjunto de datos de la figura 3.6 sería imprudente aceptar el supuesto de normalidad.

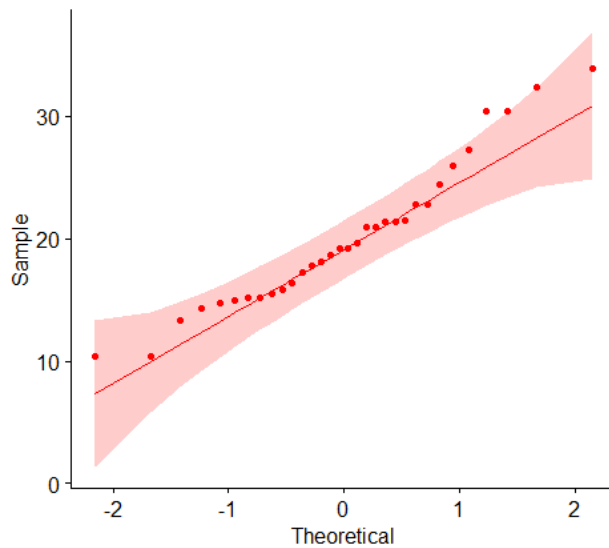


Figura 3.6: gráfico cuantil-cuantil.

Script 3.5: creación de un gráfico cuantil-cuantil.

```
1 library(ggpubr)
2
3 # Cargar datos.
4 datos <- read.csv2("C:/Inferencia/Mtcars.csv", stringsAsFactors = TRUE,
5                   row.names = 1)
6
7 # Gráfico Q-Q para la variable Rendimiento.
8 g <- ggqqplot(datos,
9               x = "Rendimiento",
10              color = "red")
11
12 print(g)
```

### 3.2.2 Distribución Z

Al trabajar con distribuciones, especialmente las simétricas, a menudo usaremos **técnicas de estandarización** para determinar qué tan usual o inusual es un determinado valor en una escala única. Así, para la distribución normal usamos como estandarización la **distribución Z** o **distribución normal estándar**, que no es más que una distribución normal centrada en 0 y con desviación estándar 1, que podemos obtener de manera bastante sencilla como muestra la ecuación 3.7.

$$Z = \frac{x - \mu}{\sigma} \quad (3.7)$$

Al aplicar la ecuación 3.7 a una observación  $x$  en una distribución normal obtenemos, entonces, su **valor**



$z$ , que determina cuán por encima o por debajo de la media (en términos de la desviación estándar) se encuentra dicha observación  $x$ . Así, observaciones cuyos valores  $z$  sean negativos estarán por debajo de la media. Análogamente, un valor  $Z$  positivo indica que la observación está por sobre la media. Mientras mayor sea el valor absoluto de su valor  $z$  ( $|z|$ ), más inusual será la observación.

### 3.2.3 Distribución chi-cuadrado

También llamada **ji-cuadrado** o  $\chi^2$ , la distribución **chi-cuadrado** (Devore, 2008) se usa para caracterizar valores siempre positivos y habitualmente desviados a la derecha. El único parámetro de esta distribución corresponde a los **grados de libertad**, usualmente representada por la letra griega  $\nu$ , que son una estimación de la cantidad de observaciones empleadas para calcular un estimador. Otra forma de entender esta idea es como la cantidad de valores que pueden cambiar libremente en un conjunto de datos. Como ejemplo, supongamos que necesitamos una muestra de tres elementos cuya media sea 10. Una vez escogidos los primeros dos, solo queda una posibilidad para el tercero de modo que se cumpla con la media deseada. Así, solo los dos primeros valores pueden cambiar libremente, por lo que se tienen dos grados de libertad.

Esta distribución está relacionada con la ya conocida distribución  $Z$ , pues si sumamos los cuadrados de  $k$  variables aleatorias independientes que siguen una distribución  $Z$ , dicha suma sigue una distribución  $\chi^2$  con  $k$  grados de libertad:

$$\sum_{i=1}^k Z_i^2 \sim \chi^2(\nu = k) \quad (3.8)$$

La media de la distribución  $\chi^2$  es  $\mu = \nu$ , y su desviación estándar,  $\sigma = 2\nu$ .

Las funciones de R para esta distribución, similares a las descritas para la distribución normal, son:

- `dchisq(x, df)`.
- `pchisq(q, df, lower.tail)`.
- `qchisq(p, df, lower.tail)`.
- `rchisq(n, df)`.

Donde:

- `x`, `q` son vectores de cuantiles (enteros no negativos).
- `p` es un vector de probabilidades.
- `n` es la cantidad de observaciones.
- `df` son los grados de libertad.
- `lower.tail` es análogo al de la función `pnorm`.

La figura 3.7 muestra un ejemplo de la distribución  $\chi^2$ .

### 3.2.4 Distribución t de Student

Ampliamente empleada cuando se trabaja con muestras pequeñas, la **distribución t de Student**, o simplemente **distribución t**, tiene, al igual que la distribución  $\chi^2$ , los grados de libertad como único parámetro. A

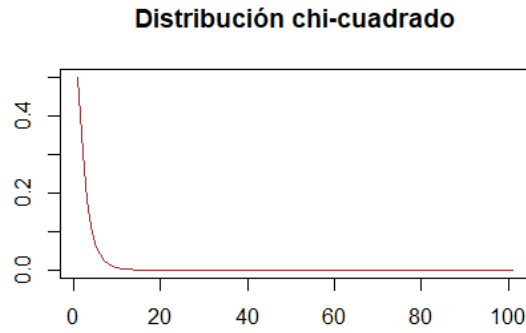


Figura 3.7: ejemplo de distribución  $\chi^2$  con 2 grados de libertad.

medida que los grados de libertad aumentan, esta distribución se asemeja cada vez más a la normal, aunque sus colas son más gruesas, como ilustra la figura 3.8.

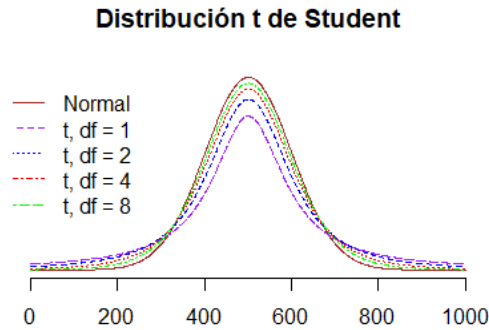


Figura 3.8: ejemplo de distribuciones t.

La distribución t se encuentra relacionada con las distribuciones vistas anteriormente de acuerdo a la ecuación 3.9, donde  $Z$  es una distribución normal estándar y  $\chi^2(\nu)$  es una distribución  $\chi^2$  con  $\nu$  grados de libertad.

$$Z \sqrt{\frac{\nu}{\chi^2(\nu)}} \sim t(\nu) \quad (3.9)$$

La media de la distribución t, para  $\nu > 1$ , es  $\mu = 0$ . Su desviación estándar, para  $\nu > 2$ , está dada por la ecuación 3.10.

$$\sigma = \sqrt{\frac{\nu}{\nu - 2}} \quad (3.10)$$

Las funciones de R para esta distribución, cuyos argumentos son análogos a los que hemos visto para las distribuciones anteriores, son:

- `dt(x, df)`.
- `pt(q, df, lower.tail)`.

- `qt(p, df, lower.tail).`
- `rt(n, df).`

### 3.2.5 Distribución F

Otra distribución que usaremos a lo largo de este libro es la **distribución F**, ampliamente usada para comparar varianzas. La distribución F se relaciona con las anteriores de acuerdo a la ecuación 3.11, donde  $\chi_1^2(\nu_1)$  y  $\chi_2^2(\nu_2)$  son dos distribuciones  $\chi^2$  con  $\nu_1$  y  $\nu_2$  grados de libertad, respectivamente. Un ejemplo de una distribución F se puede encontrar en la figura 3.9.

$$\frac{\frac{X_1^2(\nu_1)}{\nu_1}}{\frac{X_2^2(\nu_2)}{\nu_2}} \sim F(\nu_1, \nu_2) \quad (3.11)$$

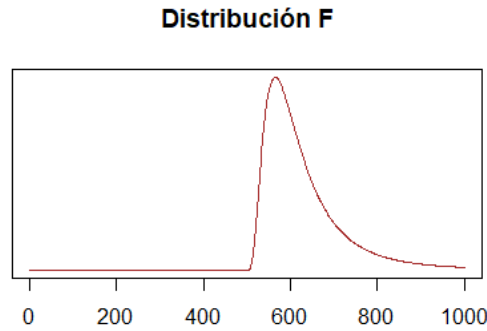


Figura 3.9: ejemplo de una distribución F.

Para  $\nu_2 > 2$ , la media de esta distribución está dada por la ecuación 3.12, y la desviación estándar corresponde a la ecuación 3.13 para  $\nu_2 > 4$ .

$$\mu = \frac{\nu_2}{\nu_2 - 2} \quad (3.12)$$

$$\sigma = \sqrt{\frac{2\nu_2^2(\nu_1 + \nu_2 - 2)}{\nu_1(\nu_2 - 2)^2(\nu_2 - 4)}} \quad (3.13)$$

Las funciones de R para esta distribución son:

- `df(x, df1, df2).`
- `pf(q, df1, df2, lower.tail).`
- `qf(p, df1, df2, lower.tail).`
- `rf(n, df1, df2).`

Donde `df1` como `df2` corresponden a grados de libertad y los argumentos restantes son los mismos que ya hemos visto anteriormente.

### 3.3 DISTRIBUCIONES DISCRETAS

Al igual que con las variables aleatorias continuas, también existen diversas distribuciones discretas de uso frecuente en estadística.

#### 3.3.1 Distribución de Bernoulli

Una **variable aleatoria de Bernoulli** es aquella en que cada intento individual tiene solo dos resultados posibles: “éxito”, que ocurre con una probabilidad  $p$  y se representa habitualmente con un 1, y “fracaso”, que ocurre con probabilidad  $q = 1 - p$  y suele representarse por un 0. La selección de qué resultado se considera como éxito o fracaso suele ser arbitraria. Para ilustrar esta idea, si dos personas lanzan una moneda al aire para sortear al ganador, cada una de ellas considerará una cara diferente de la moneda como un éxito.

Otro ejemplo que nos puede ayudar es el de lanzar varios dados de 20 caras, donde el éxito corresponda a obtener un 20 como resultado. Cada uno de ellos tiene una **probabilidad de éxito** (obtener 20)  $p = 0.05$  y una **probabilidad de fracaso** (obtener otro valor)  $q = 1 - p = 0.95$ . Los lanzamientos de los dados son **independientes**, pues un dado no afecta a los demás.

Definimos la **proporción de la muestra** para una distribución de Bernoulli,  $\hat{p}$ , como la cantidad de éxitos dividida por la cantidad de intentos. Mientras mayor sea la cantidad de intentos, más cercano será el valor de  $\hat{p}$  a la probabilidad real de éxito  $p$ .

Al igual que la distribución normal, la distribución de Bernoulli puede resumirse expresando su media ( $\mu = p$ ) y su desviación estándar. Esta última está dada por la ecuación 3.14.

$$\sigma = \sqrt{p(1-p)} \quad (3.14)$$

El paquete `extraDistr` de R ofrece 4 funciones, similares a las ya conocidas, para la distribución de Bernoulli:

- `dbern(x, prob)`.
- `pbern(q, prob, lower.tail)`.
- `qbern(p, prob, lower.tail)`.
- `rbern(n, prob)`.

#### 3.3.2 Distribución geométrica

La **distribución geométrica** describe la cantidad de intentos que debemos realizar hasta obtener un éxito para variables de Bernoulli **independientes e idénticamente distribuidas**, es decir, que no se afectan unas a otras y cada una con igual probabilidad de éxito.

La probabilidad de obtener un éxito al  $n$ -ésimo intento está dada por la ecuación 3.15, donde podemos ver que las probabilidades en esta distribución decrecen exponencialmente rápido, como ilustra la figura 3.10. La media y la desviación estándar de la distribución geométrica están dadas, respectivamente, por las ecuaciones 3.16 y 3.17.

$$\Pr(\text{primer éxito al } n\text{-ésimo intento}) = (1-p)^{n-1}p \quad (3.15)$$

$$\mu = \frac{1}{p} \quad (3.16)$$

$$\sigma = \sqrt{\frac{1-p}{p^2}} \quad (3.17)$$

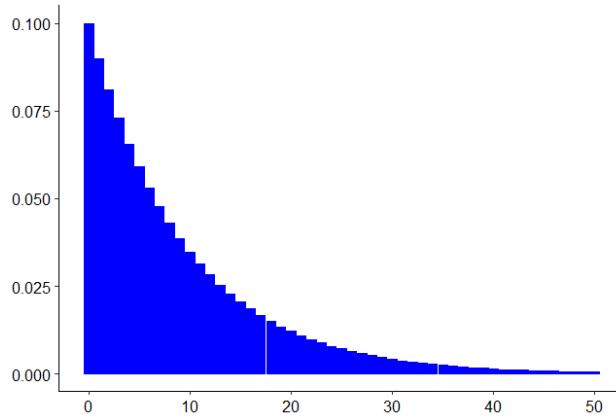


Figura 3.10: distribución geométrica para obtener un valor específico lanzando un dado de 20 caras.

Para entender mejor la utilidad de la distribución geométrica, consideremos la pregunta: ¿cuántas veces tenemos que lanzar un dado de 20 caras para obtener un 1? Anteriormente vimos que la probabilidad de éxito en este caso es  $p = 0.05$ . El valor esperado, representado por la media, sería en este caso el que se presenta en la ecuación 3.18.

$$\mu = \frac{1}{p} = \frac{1}{0.05} = 20 \quad (3.18)$$

Una vez más, R ofrece funciones similares a las presentadas anteriormente para trabajar con distribuciones geométricas:

- `dgeom(x, prob)`.
- `pgeom(q, prob, lower.tail)`.
- `qgeom(p, prob, lower.tail)`.
- `rbern(n, prob)`.

### 3.3.3 Distribución binomial

A diferencia de la distribución geométrica, la **distribución binomial** describe la probabilidad de tener exactamente  $k$  éxitos en  $n$  intentos independientes de Bernoulli con probabilidad de éxito  $p$ , cuya función de probabilidad está dada por la ecuación 3.19, donde:

- $\binom{n}{k} = \frac{n!}{k!(n-k)!}$  corresponde a la cantidad de formas de obtener  $k$  éxitos en un total de  $n$  intentos.
- $p^k(1-p)^{n-k}$  es la probabilidad de tener un único éxito en solo una de las  $\binom{n}{k}$  maneras posibles.

$$f(k, n, p) = \Pr(k; n, p) = \Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k} \quad (3.19)$$

La media y la desviación estándar de la distribución binomial están dadas por las ecuaciones 3.20 y 3.21, respectivamente. Un ejemplo de esta distribución se presenta en la figura 3.11

$$\mu = np \quad (3.20)$$

$$\sigma = \sqrt{np(1 - p)} \quad (3.21)$$

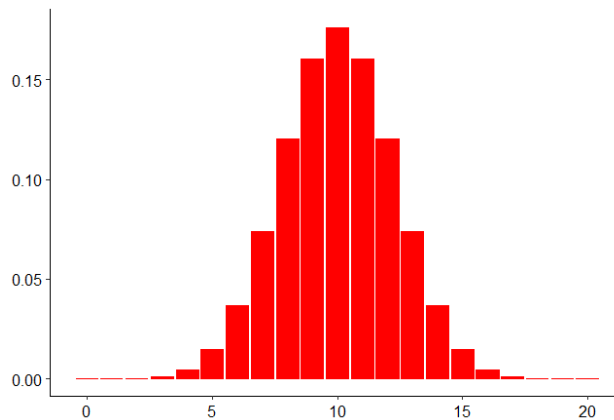


Figura 3.11: distribución binomial con  $\mu = 400$  y  $\sigma = 15.4019$ .

Antes de decidir usar la distribución binomial, tenemos que verificar cuatro condiciones:

1. Los intentos son independientes.
2. La cantidad de intentos ( $n$ ) es fija.
3. El resultado de cada intento puede ser clasificado como éxito o fracaso.
4. La probabilidad de éxito ( $p$ ) es la misma para cada intento.

Las funciones que ofrece R para trabajar con esta distribución son:

- `dbinom(x, size, prob)`.
- `pbinom(x, size, prob)`.
- `qbinom(p, size, prob)`.
- `rbinom(n, size, prob)`.

Donde:

- **x** es un vector numérico.
- **p** es un vector de probabilidades.
- **n** es la cantidad de observaciones.
- **size** corresponde al número de intentos.
- **prob** es la probabilidad de éxito de cada intento.

En la figura 3.11 podemos observar que, en cierto modo, la distribución binomial se asemeja a la distribución normal: ambas son simétricas, aunque la distribución binomial no tiene la forma de campana de la distribución normal. Esta similitud ofrece una importante ventaja, pues en ocasiones es posible usar la distribución normal para estimar probabilidades binomiales, evitando así el uso de la compleja fórmula de la ecuación 3.19. Formalmente, esta aproximación es válida cuando el tamaño de la muestra,  $n$ , es lo suficientemente grande para que tanto  $np$  como  $n(1 - p)$  sean mayores o iguales que 10. En este caso, los parámetros de la distribución normal aproximada son los mismos de la distribución binomial (ecuaciones 3.20 y 3.21).

### 3.3.4 Distribución binomial negativa

La **distribución binomial negativa** es algo más general que la binomial, pues describe la probabilidad de encontrar el  $k$ -ésimo éxito al  $n$ -ésimo intento. Como señalan Diez y col. (2017, p. 155), “en el caso binomial, en general se tiene una cantidad fija de intentos y se considera la cantidad de éxitos. En el caso binomial negativo, se examina cuántos intentos se necesitan para observar una cantidad fija de éxitos y se requiere que la última observación sea un éxito”<sup>3</sup>.

Como adelanta la comparación anterior, antes de decidir usar la distribución binomial negativa tenemos que verificar cuatro condiciones:

1. Los intentos son independientes.
2. El resultado de cada intento puede ser clasificado como éxito o fracaso.
3. La probabilidad de éxito ( $p$ ) es la misma para cada intento.
4. El último intento debe ser un éxito.

La función de probabilidad para esta distribución, ejemplificada en la figura 3.12, está dada por la ecuación 3.22. La varianza y la desviación estándar están dadas por las ecuaciones 3.23 y 3.24 (Devore, 2008, p. 120).

$$\Pr(k\text{-ésimo éxito al } n\text{-ésimo intento}) = \binom{n-1}{k-1} p^k (1-p)^{n-k} \quad (3.22)$$

$$\mu = \frac{k(1-p)}{p} \quad (3.23)$$

$$\sigma = \sqrt{\frac{k(1-p)}{p^2}} \quad (3.24)$$

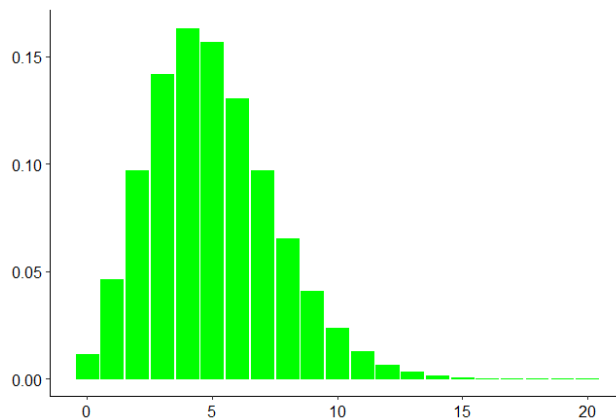


Figura 3.12: ejemplo de distribución binomial negativa.

Nuevamente, R dispone de cuatro funciones que permiten trabajar con esta distribución:

- `dnbinom(x, size, prob)`.
- `pnbinom(q, size, prob, lower.tail)`.
- `qnbinom(p, size, prob, lower.tail)`.
- `rnbinom(n, size, prob)`.

Donde:

---

<sup>3</sup>Traducción libre de los autores.

- `x`, `q` son vectores de cuantiles (enteros no negativos).
- `p` es un vector de probabilidades.
- `n` es la cantidad de observaciones.
- `size` corresponde al número (no negativo) de intentos.
- `prob` es la probabilidad de éxito de cada intento.
- `lower.tail` es análogo al de la función `pnorm`.

### 3.3.5 Distribución de Poisson

Útil para estimar la cantidad de eventos en una población grande en un lapso de tiempo dado, por ejemplo, la cantidad de contagios de influenza entre los habitantes de Santiago en una semana, la **distribución de Poisson** (figura 3.13) tiene una función de probabilidad definida por la ecuación 3.25, donde  $\lambda$  es la tasa o cantidad de eventos que se espera observar en un lapso de tiempo dado y  $k$  puede tomar cualquier valor entero no negativo. La media de esta distribución está dada por  $\lambda$  y la desviación estándar, por  $\sqrt{\lambda}$ .

$$\Pr(\text{observar } k \text{ eventos}) = \frac{\lambda^k e^{-\lambda}}{k!} \quad (3.25)$$

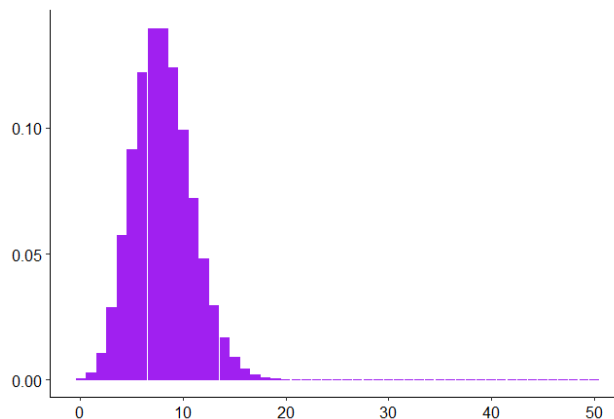


Figura 3.13: ejemplo de distribución de Poisson.

Las funciones de R para esta distribución son:

- `dpois(x, lambda)`.
- `ppois(q, lambda, lower.tail)`.
- `qppois(p, lambda, lower.tail)`.
- `rpois(n, lambda)`.

Donde:

- `x`, `q` son vectores de cuantiles (enteros no negativos).
- `p` es un vector de probabilidades.
- `n` es la cantidad de observaciones.
- `lambda` es un vector no negativo de medias.
- `lower.tail` es análogo al de la función `pnorm`.



### 3.4 EJERCICIOS PROPUESTOS

1. Da un ejemplo de variable aleatoria (novedosa) que puedas observar en tus compañeros y que tenga una función de densidad de probabilidad discreta.
2. Para la variable anterior, ¿cuál sería el valor esperado? ¿Cuál sería la varianza? ¿Cómo te imaginas su función de densidad de probabilidad?
3. Lista tres nombres distintos con que también se llama a las funciones de densidad de probabilidad.
4. Si una variable aleatoria tiene una función de densidad de probabilidad con media igual a 30 y desviación estándar de 3, ¿por qué podría ocurrir que la probabilidad de que la variable tome el valor 30 sea nula, es decir,  $P(X = 30) = 0$ ?
5. Según el Reporte Mensual de Empleo, las siguientes son las estadísticas (media  $\pm$  desviación estándar) para las seis variables relevantes que se han estudiado en los últimos cinco años:
  - a) Número de personas despedidas:  $64.675 \pm 8.321$ .
  - b) Número de personas renunciadas:  $118.543 \pm 17.936$ .
  - c) Número de personas jubiladas:  $97.092 \pm 11.147$ .
  - d) Número de empleos creados:  $24.715 \pm 10.832$ .
  - e) Número de personas contratadas:  $301.345 \pm 27.261$ .
  - f) Número de personas entrando a la fuerza de trabajo:  $26.444 \pm 29.440$ .Con esta información, calcula la media y la desviación estándar de:
  - a) Caída neta del empleo:  $(d) - (a) - (b) - (c)$ .
  - b) Subida neta del empleo:  $(e) - (a) - (b) - (c)$ .
  - c) Caída neta del desempleo:  $(a) + (b) + (e) + (f)$ .
  - d) Vacancia del empleo:  $(d) - (e)$ .
6. ¿Qué significa que cierto valor de una variable aleatoria, usualmente con distribución normal, tenga valor  $Z = 1, 5$ ?
7. Según la regla empírica, ¿entre qué estaturas se podría encontrar al 95 % de los estudiantes varones del Departamento de Ingeniería Informática de la Universidad de Santiago de Chile, si esta variable sigue una distribución  $N(\mu = 171, \sigma = 3)$ ?
8. ¿Qué información entrega un Gráfico Q-Q? ¿Para qué se usa?
9. La probabilidad de que un estudiante universitario chileno seleccionado al azar sea VIH positivo es 0,013. ¿Cuáles serían la media y la desviación estándar de esta variable?
10. En promedio, ¿a cuántos estudiantes universitarios se debería revisar hasta encontrar a un VIH positivo?
11. Si el Departamento de Salud de una Universidad chilena controla a 50 estudiantes por día durante una semana de clases (lunes a viernes), ¿cuál sería el número promedio de VIH positivos detectados cada día? ¿Con qué varianza?
12. Si la Universidad del ejercicio anterior dispone de 10 paquetes de tratamiento de VIH para estudiantes, ¿cómo podría saber a cuántos estudiantes debería examinar para poder asignarlos (suponiendo que todo estudiante VIH positivo acepta el tratamiento)?
13. Muestra un ejemplo novedoso de una variable aleatoria relacionada que podría seguir una distribución de Poisson.



## CAPÍTULO 4. FUNDAMENTOS PARA LA INFERENCIA

En el capítulo 1 se definen los conceptos de población, entendido como todo el conjunto de interés, y muestra, que es un subconjunto de la población. También se introducen las nociones de parámetro, correspondiente a un valor que resume la población (por ejemplo la media de la población,  $\mu$ ), y de estadístico, como valor que resume una muestra (por ejemplo, la media muestral,  $\bar{x}$ ). La **inferencia estadística** tiene por objeto entender cuán cerca está el estadístico del parámetro real de la población. En este capítulo conoceremos los principios necesarios para la inferencia estadística, con base en Diez y col. (2017, pp. 168-202) y Field y col. (2012, pp. 40-47).

### 4.1 ESTIMADORES PUNTUALES

Como ya dijimos, los parámetros y los estadísticos son valores que resumen, respectivamente, una población y una muestra. En consecuencia, podemos decir que un estadístico corresponde a un **estimador puntual** de un parámetro. El valor de un estimador puntual cambia dependiendo de la muestra que usemos para obtenerlo. Así, por más que su valor se acerque al parámetro de la población, difícilmente será igual a este último. Sin embargo, el estimador tiende a mejorar a medida que aumentamos el tamaño de la muestra, por efecto de la **ley de los grandes números**. Para ilustrar este fenómeno, consideremos la **media móvil**, que es una secuencia de medias muestrales en que cada una de ellas toma un elemento más de la población que su antecesora. La figura 4.1, elaborada con el script 4.1, ejemplifica este fenómeno.

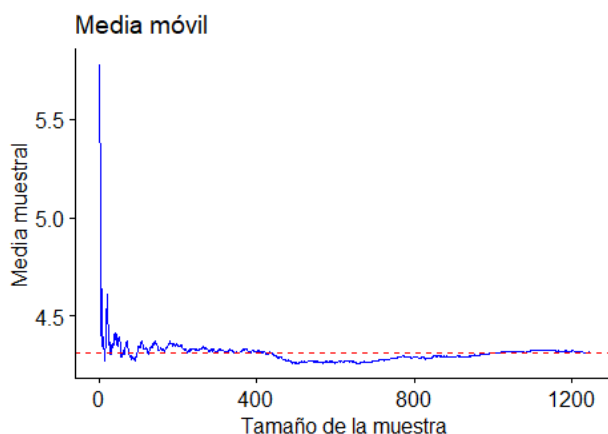


Figura 4.1: medias obtenidas al agregar a la muestra un elemento cada vez.

Script 4.1: representación gráfica de la media móvil.

```
1 library(ggpubr)
2
3 # Establecer la semilla para generar números aleatorios.
4 set.seed(9437)
5
6 # Generar aleatoriamente una población de tamaño 1500
7 # (en este caso, con una distribución cercana a la normal).
```

```

8 poblacion <- rnorm(n = 1500, mean = 4.32, sd = 0.98)
9
10 # Calcular la media de la población.
11 media_poblacion <- mean(poblacion)
12 cat("Media de la población:", media_poblacion, "\n")
13
14 # Tomar una muestra de tamaño 1250.
15 tamano_muestra <- 1250
16 muestra <- sample(poblacion, tamano_muestra)
17
18 # Calcular las medias acumuladas (es decir, con muestras de
19 # 1, 2, 3, ... elementos).
20 n <- seq(along = muestra)
21 media <- cumsum(muestra) / n
22
23 # Crear una matriz de datos con los tamaños y las medias muestrales.
24 datos <- data.frame(n, media)
25
26 # Graficar las medias muestrales.
27 g <- ggline(data = datos,
28             x = "n",
29             y = "media",
30             plot_type = "l",
31             color = "blue",
32             main = "Media móvil",
33             xlab = "Tamaño de la muestra",
34             ylab = "Media muestral")
35
36 # Añadir al gráfico una recta con la media de la población.
37 g <- g + geom_hline(aes(yintercept = media_poblacion),
38                    color = "red", linetype = 2)
39
40 print(g)

```

Para determinar qué tan adecuado es un estimador, necesitamos saber cuánto cambia de una muestra a otra. Si esta variabilidad es pequeña, es muy probable que la estimación sea buena. Podemos estudiar la variabilidad de la muestra con ayuda de la **distribución muestral**, que representa la distribución de estimadores puntuales obtenidos con **todas** las diferentes muestras de igual tamaño de una misma población. La figura 4.2 (construida con el script 4.2) representa las medias para diferentes muestras de una población, aunque solo una selección aleatoria de todas las posibles muestras, incluyendo además una línea vertical roja que señala la media de la población. Podemos destacar que las medias muestrales tienden a aglutinarse en torno a la media poblacional, pues de acuerdo al **teorema del límite central**, la distribución de  $\bar{x}$  se aproxima a la normalidad. Esta aproximación mejora a medida que aumenta el tamaño de la muestra.

Script 4.2: distribución de la media muestral.

```

1 library(ggpubr)
2
3 # Establecer la semilla para generar números aleatorios.
4 set.seed(94)
5
6 # Generar aleatoriamente una población de tamaño 1500
7 # (en este caso, con una distribución cercana a la normal).
8 poblacion <- rnorm(n = 1500, mean = 4.32, sd = 0.98)
9
10 # Calcular la media de la población.
11 media_poblacion <- mean(poblacion)
12 cat("Media de la población:", media_poblacion, "\n")
13

```

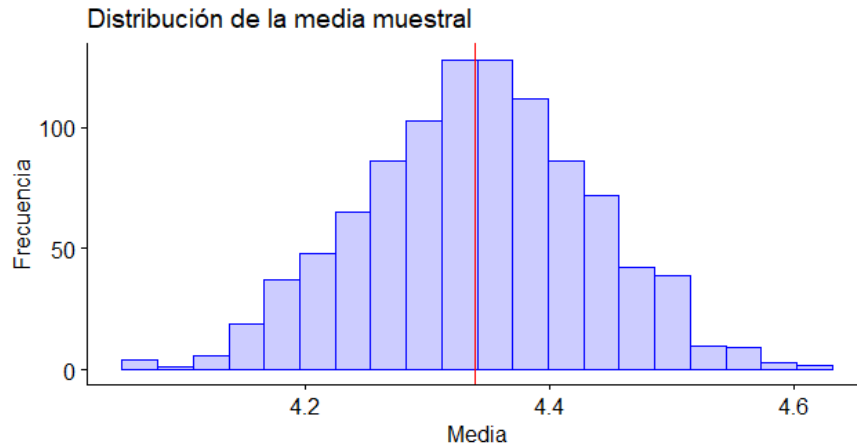


Figura 4.2: distribución muestral de la media para muestras con 100 observaciones.

```

14 # Tomar 1000 muestras de tamaño 100. Quedan almacenadas
15 # como una matriz donde cada columna es una muestra.
16 tamaño_muestra <- 100
17 repeticiones <- 1000
18
19 muestras <- replicate(repeticiones,
20                       sample(poblacion, tamaño_muestra))
21
22 # Calcular medias muestrales y almacenar los resultados
23 # en forma de data frame.
24 medias <- colMeans(muestras)
25 medias <- as.data.frame(medias)
26
27 # Construir un histograma de las medias muestrales.
28 g <- gghistogram(data = medias,
29                  x = "medias",
30                  bins = 20,
31                  title = "Distribución de la media muestral",
32                  xlab = "Media",
33                  ylab = "Frecuencia",
34                  color = "blue",
35                  fill = "blue",
36                  alpha = 0.2)
37
38 # Agregar línea vertical con la media de la población.
39 g <- g + geom_vline(aes(xintercept = media_poblacion),
40                    color = "red", linetype = 1)
41
42 print(g)

```

## 4.2 MODELOS ESTADÍSTICOS

Ahora que hemos conocido más conceptos, podemos definir con precisión qué es un **modelo estadístico**. En el capítulo 1 dijimos que un modelo es simplemente una representación y que los modelos estadísticos pueden

emplearse para diversos propósitos:

- Describir o resumir datos.
- Clasificar objetos o predecir resultados.
- Anticipar los resultados de intervenciones (en ocasiones).

Más formalmente, un modelo estadístico es una descripción de un **proceso probabilístico** con **parámetros desconocidos** que deben ser **estimados** en base a **suposiciones** y un conjunto de datos **observados**. En general, tiene la forma dada en la ecuación 4.1:

$$y_i = (\text{modelo}) + \varepsilon_i \quad (4.1)$$

Donde:

- $y_i$  es el  $i$ -ésimo valor observado de la variable respuesta  $Y$  (también llamada variable de salida o variable dependiente).
- modelo es el resultado de una función determinista basada en un conjunto de argumentos.
- $\varepsilon_i$  es el error, correspondiente a la **variación natural**, y no a una equivocación, existente entre los valores observados y los valores pronosticados por el modelo. También recibe los nombres de variación no sistemática, variación aleatoria, residuos o incluso, residuales.

El error  $\varepsilon_i$  en la ecuación 4.1 se relaciona entonces con la calidad del modelo. Mientras menor sea el error, mejor será el modelo. Por el contrario, un error grande es señal de un modelo fallido, que no describe bien los datos, no ayuda a predecirlos bien, o no ayuda a su correcta clasificación.

La media y la proporción, y cualquier estadístico en general, son, en sí mismos, modelos estadísticos, aunque bastante simples.

## 4.3 ERROR ESTÁNDAR

En el capítulo 2 conocimos la desviación estándar como medida que estima la distancia de las observaciones respecto de la media. El **error estándar**, denotado usualmente por  $SE_{\hat{\theta}}$  o  $\sigma_{\hat{\theta}}$ , corresponde a la desviación estándar de la distribución de un estimador muestral  $\hat{\theta}$  de un parámetro  $\theta$ . Por ejemplo, el error estándar de la media, es decir la desviación estándar de la distribución de las medias de todas las posibles muestras de  $n$  observaciones independientes, se calcula de acuerdo a la ecuación 4.2.

$$SE_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \quad (4.2)$$

Donde  $\sigma$  es la desviación estándar de la población y  $n$  corresponde al tamaño de la muestra. En esta ecuación queda en evidencia que el error estándar de la media disminuye a medida que el tamaño de la muestra aumenta. Un método confiable que podemos usar para asegurar que las observaciones sean independientes es realizar un muestreo aleatorio simple<sup>1</sup> que abarque menos del 10 % de la población.

Volviendo a la ecuación para calcular el error estándar de la media muestral (ecuación 4.2), ¡debemos tener cuidado antes de usarla! Ya hemos mencionado antes que la distribución de las medias muestrales tiende a ser cercana a la normal, por lo que en dicho caso es posible usar el **modelo normal**, sustentado en el teorema del límite central. Las condiciones que deben cumplirse para usar este modelo y que, en consecuencia, el error estándar sea preciso, son:

1. Las observaciones de la muestra son independientes.

---

<sup>1</sup>Es decir, una muestra en que todos los elementos de la población tengan igual probabilidad de ser escogidos. Las técnicas de muestreo se abordan con más detalle en capítulos posteriores.

2. La muestra es grande (en general  $n \geq 30$ ).
3. La distribución de la muestra no es significativamente asimétrica. Esto último suele además relacionarse con la presencia de valores atípicos. Mientras mayor sea el tamaño de la muestra, más se puede relajar esta condición.

Si no se cumplen las condiciones anteriores, debemos considerar otras opciones: para muestras pequeñas, se deben considerar métodos alternativos, y si la distribución de la muestra presenta una asimetría significativa, entonces tendremos que incrementar el tamaño de la muestra para compensar el efecto de la desviación.

## 4.4 INTERVALOS DE CONFIANZA

Hasta ahora sabemos que un estimador puntual es un único valor (obtenido a partir de una muestra) que, como su nombre indica, estima un parámetro de la población. Por ende, dicho valor rara vez es exacto. En consecuencia, lo lógico sería establecer un rango de valores plausibles para el parámetro estimado, que llamaremos **intervalo de confianza**, y que se construye en torno al estimador puntual. Dado que el error estándar representa la desviación estándar asociada al estimador, tiene sentido que lo usemos como guía en este proceso.

Recordemos que en el capítulo 3 vimos una regla empírica para la distribución normal (figura 3.5), la cual señala que (para distribuciones normales) alrededor de 95 % de las veces el estimador puntual se encontrará en un rango de 2 errores estándar del parámetro. Es decir, al considerar un intervalo de confianza de dos errores estándar (4.3), tendremos 95 % de **confianza** de haber capturado el parámetro real.

$$\bar{x} \pm 2 \cdot SE_{\bar{x}} \quad (4.3)$$

Podemos generalizar la ecuación 4.3 para calcular el intervalo de confianza para la media con cualquier **nivel de confianza** como muestra la ecuación 4.4.

$$\bar{x} \pm z^* \cdot SE_{\bar{x}} \quad (4.4)$$

El término  $z^*$  en la ecuación 4.4 corresponde, usualmente, al valor  $z$  tal que el área bajo la curva normal estándar comprendida entre  $-z^*$  y  $z^*$  es igual al nivel de confianza deseado. La expresión  $z^* \cdot SE$  recibe el nombre de **margen de error**.

Tomemos como ejemplo un **nivel de confianza** (que, por razones que veremos en la sección siguiente, denotaremos por  $1 - \alpha$ ) de 90 % (es decir,  $1 - \alpha = 0,9$ ). Eso significa, entonces, que nuestro intervalo de confianza excluye el 5 % del área correspondiente a la cola inferior (es decir, el percentil con valor 0,05) e igual porcentaje del área correspondiente a la cola superior (que, como la distribución  $Z$  es simétrica, es igual al área anterior). Puesto que conocemos el percentil,  $(1 - \alpha)/2 = 0,05$ , en R podemos usar la llamada `qnorm(0.05, mean = 0, sd = 1, lower.tail = FALSE)` y obtenemos  $z^* = 1,64$ . Es importante indicar que en esta llamada estamos en realidad trabajando con la cola superior para que  $z^*$  sea positivo. Si hacemos la llamada para la cola inferior, obtenemos  $z^* = -1,64$ .

Es importante destacar que, una vez más, debemos ser cuidadosos al interpretar un intervalo de confianza del  $x$  % ( $x = 1 - \alpha$ ). Su significado es, sencillamente, “se tiene  $x$  % de certeza de que el parámetro de la población se encuentra entre...” (Diez y col., 2017, p. 180), es decir, que, en promedio,  $x$  % de los intervalos de confianza que se construyan en torno a un estadístico, con muestras de un tamaño fijo, capturarán el verdadero valor del parámetro. Esto **no es equivalente** a decir que el valor del parámetro tiene una “probabilidad de  $x$  %” de estar entre los valores del intervalo calculado, lo que sería incorrecto. Por otra parte, los intervalos de confianza no dicen nada acerca de observaciones individuales, sino que solo hablan del parámetro en cuestión.

## 4.5 PRUEBAS DE HIPÓTESIS

Supongamos que un banco ha desarrollado un nuevo sistema computacional para gestionar sus transacciones. El nuevo sistema ( $N$ ) se ha puesto a prueba durante un mes, funcionando (con iguales condiciones de hardware) en paralelo con el sistema antiguo ( $A$ ) y el banco ha llevado un registro del tiempo que tarda cada sistema en efectuar cada transacción. El gerente ha determinado que autorizará la migración al nuevo sistema únicamente si este es más rápido que el antiguo para procesar las transacciones. Se sabe que el sistema antiguo tarda en promedio  $\mu_A = 530$  milisegundos en procesar una transacción. Para el sistema nuevo se han registrado  $n = 1.600$  transacciones, realizadas en un tiempo promedio de  $\bar{x}_N = 527,9$  [ms] con desviación estándar  $s_N = 48$  [ms].

Una primera aproximación para tomar la decisión puede ser investigar si existe diferencia en los tiempos de ejecución de ambos sistemas, lo que puede expresarse en torno a dos **hipótesis** (palabra que la Real Academia Española (2014) define como “Suposición de algo posible o imposible para sacar de ello una consecuencia”) que compiten entre sí:

$H_0$ : El nuevo sistema, en promedio, tarda lo mismo que el antiguo en procesar las transacciones, es decir:  
 $\mu_N = \mu_A$ .

$H_A$ : Los sistemas requieren, en promedio, cantidades de tiempo diferentes para procesar las transacciones, es decir:  $\mu_N \neq \mu_A$

La primera hipótesis,  $H_0$ , recibe el nombre de **hipótesis nula** y suele representar una postura escéptica, es decir, que no hay cambios, por lo que **la hipótesis nula siempre se formula como una igualdad!**. La segunda ( $H_A$ ), llamada **hipótesis alternativa**, representa en cambio una nueva perspectiva. Esta primera aproximación corresponde a una **prueba bilateral** o de dos colas, pues la diferencia puede ser en ambos sentidos:  $H_0$  no parece correcta si  $\mu_N < \mu_A$  o si  $\mu_N > \mu_A$ .

Como en este caso conocemos el valor de  $\mu_A = 530$  [ms], también podríamos escribir la formulación matemática de las hipótesis de la siguiente manera:

$H_0$ :  $\mu_N = 530$

$H_A$ :  $\mu_N \neq 530$

En este planteamiento, “530” recibe el nombre de **valor nulo**, pues representa el valor del parámetro cuando se cumple la hipótesis nula.

Una aproximación más cercana al problema descrito puede ser investigar si el nuevo sistema es efectivamente **más rápido** que el antiguo. En este caso, se habla de una **prueba unilateral** o de una cola, pues solo interesa saber si el tiempo promedio empleado por el nuevo sistema es menor que el empleado por el sistema antiguo. Las hipótesis, en este caso, serían:

$H_0$ : El nuevo sistema tarda, en promedio, lo mismo que el antiguo en procesar las transacciones, es decir:  
 $\mu_N = \mu_A$ .

$H_A$ : El nuevo sistema tarda, en promedio, menos que el antiguo en procesar las transacciones, es decir:  
 $\mu_N < \mu_A$

Obviamente en otros casos podría interesar solamente si valor alternativo es mayor que el valor nulo.

Teniendo las hipótesis planteadas, sigue decidir si la hipótesis nula parece o no plausible a través de una **prueba de hipótesis**. El marco para la prueba de hipótesis es **escéptico**: no se rechaza la hipótesis nula a menos que haya suficiente evidencia para rechazarla en favor de la hipótesis alternativa. Esta idea es muy parecida a la expresada en la expresión de uso común “se presume inocente mientras no se demuestre lo contrario”. Sin embargo, el que no se logre rechazar  $H_0$  **no significa aceptarla** como verdadera o como correcta sin más. Por eso se usa un lenguaje bastante peculiar, señalando que *se falla al rechazar  $H_0$*  o bien que *se rechaza  $H_0$  en favor de  $H_A$* . Retomando la analogía con la expresión anterior, que no haya pruebas suficientes para la culpabilidad, no significa que una persona sea en verdad inocente.

Volvamos al escenario del ejemplo para la prueba de hipótesis bilateral (es decir, aquella en que solo queremos



ver si hay diferencias en el tiempo de procesamiento de transacciones entre ambos sistemas del banco). El valor de  $\bar{x}_N = 527,9$  [ms] es, en efecto, distinto de  $\mu_A = 530$  [ms]. No obstante, al ser una estimación puntual, como ya hemos estudiado, esta diferencia podría deberse simplemente a la muestra escogida, por lo que el parámetro real  $\mu_N$  podría ser igual a  $\mu_A$  [ms]. En consecuencia, resulta útil calcular el intervalo de confianza para  $\bar{x}_N$ .

Comencemos por determinar el error estándar:

$$SE_{\bar{x}} = \frac{s_N}{\sqrt{n}} = \frac{48}{\sqrt{1600}} = 1,2$$

Ahora fijemos un nivel de confianza, por ejemplo 95 %, y usemos el valor  $z^*$  correspondiente para calcular el intervalo de confianza:

$$\bar{x}_N \pm z^* \cdot SE_{\bar{x}} = 527,9 \pm 1,96 \cdot 1,2 = [525,548; 530,252]$$

Como el parámetro del sistema antiguo ( $\mu_A = 530$  [ms]) cae (a penas) dentro de este intervalo, se puede suponer que no existe una diferencia significativa entre los tiempos promedio requeridos por ambos sistemas, por lo que no se rechaza  $H_0$ . Así, tenemos un 95 % de confianza en que no existe una diferencia entre los tiempos que requieren ambos sistemas para procesar transacciones. Sin embargo, esta decisión es un tanto apresurada ya que el resultado está cerca del borde de rechazo y, en este caso, lo lógico sería investigar más (hacer crecer la muestra).

Revisemos ahora el caso planteado con hipótesis alternativa unilateral (es decir, queremos ver si el nuevo sistema es, en efecto, más rápido). Manteniendo nuestro nivel de confianza  $1 - \alpha = 0,95$ , en este caso debemos considerar los valores menores a  $\mu_A = 530$  [ms] para el cálculo de  $z^*$ . En otras palabras, el 5 % que descartamos corresponde únicamente a la cola superior. Así, nuestro valor para  $z^*$  está dado por la llamada `qnorm(0.05, mean = 0, sd = 1, lower.tail = FALSE)`, obteniéndose  $z^* = 1,64$  (aprox.) por lo que se tiene que la cota superior es:

$$\bar{x}_N - z^* \cdot SE_{\bar{x}} = 527,9 - 1,64 \cdot 1,2 = 529.874$$

Luego, el intervalo de confianza va desde “cualquier valor” bajo la media observada en la muestra hasta el valor calculado arriba, por lo que el intervalo con 95 % confianza sería:  $[-\infty; 529,874]$ .

Ahora el valor  $\mu_A = 530$  [ms] cae (apenas) fuera del intervalo y podemos decir que existe evidencia de que el nuevo sistema tarda en promedio menos tiempo que el antiguo en procesar las transacciones.

Ahora bien, siempre que se prueban hipótesis podemos cometer un error al momento de decidir si rechazar o no la hipótesis nula. Afortunadamente, la estadística ofrece herramientas para cuantificar cuán frecuentes son dichos errores. Existen cuatro posibles escenarios, los cuales se presentan en la tabla 4.1. El **error tipo I** corresponde a rechazar  $H_0$  cuando en realidad es verdadera, mientras que el **error tipo II** corresponde a no rechazarla cuando en realidad  $H_A$  es verdadera.

		Conclusión de la prueba	
		No rechazar $H_0$	Rechazar $H_0$ en favor de $H_A$
Verdad	$H_0$ verdadera	Decisión correcta	Error tipo I
	$H_A$ verdadera	Error tipo II	Decisión correcta

Tabla 4.1: posibles escenarios para una prueba de hipótesis.

Como ya hemos señalado, la prueba de hipótesis se basa en no rechazar  $H_0$  a menos que se tenga evidencia contundente. Por regla general, no se desea cometer el error de rechazar incorrectamente la hipótesis nula (error tipo I) en más de 5 % de los casos. Esto corresponde a un **nivel de significación** de 0,05, denotado por  $\alpha = 0,05$ . Si usamos un intervalo de confianza de 95 % para evaluar una prueba de hipótesis en que la

hipótesis nula es verdadera, cometeremos un error tipo I cada vez que el estimador puntual esté a 1,96 o más errores estándar del parámetro de la población. Esto puede ocurrir un 5 % de las veces (2,5 % en cada cola de la distribución para el caso bilateral). Del mismo modo, un intervalo de confianza del 99 % es equivalente a un nivel de significación  $\alpha = 0,01$ .

El intervalo de confianza es de mucha ayuda para decidir si rechazar o no  $H_0$ . No obstante, no aporta información directa acerca de cuán fuerte es la evidencia para la decisión tomada.

#### 4.5.1 Prueba formal de hipótesis con valores p

Antes de que la computación se hiciera masiva, las personas tenían dos procedimientos posibles para decidir una prueba de hipótesis. El primero es el realizado en la sección anterior, esto es, calcular el intervalo con  $(1 - \alpha) \%$  de confianza de acuerdo a los estadísticos observados en una muestra y revisar si el valor nulo cae o no dentro de este intervalo. El otro procedimiento clásico, que podemos encontrar en muchos libros y sitios en Internet, es estimar a qué valor  $z$  corresponde la media observada en la distribución normal estandarizada que define el valor nulo y el error estándar: si este estadístico  $z$  es mayor que  $z^*$ , entonces el estadístico cae en una “zona de rechazo” de  $H_0$ ; en caso contrario ( $|z| < z^*$ ), se falla en rechazar la hipótesis nula.

Si bien estos procedimientos siguen siendo útiles, su diseño respondía a la existencia de **tablas de probabilidad** en que se tabulaban probabilidades para algunos valores de percentiles de uso común, como 90 %, 95 %, 0,975 % o 0,99 %.

Con la llegada de los computadores, y en particular de entornos como R, es posible obtener probabilidades (casi) exactas para cualquier percentil. Esto hizo que un tercer método para decidir una prueba de hipótesis haya ido ganando popularidad: el uso del **valor p**, también llamado **p-valor**, que es definido por Diez y col. (2017, p. 186) como “la probabilidad de observar datos al menos tan favorables como la muestra actual para la hipótesis alternativa, si la hipótesis nula es verdadera”. De esta forma, un p-valor permite cuantificar cuán fuerte es la evidencia en contra de la hipótesis nula (y en favor de la hipótesis alternativa).

Consideremos ahora el escenario de la hipótesis unilateral del ejemplo, con un nivel de significación  $\alpha = 0,05$ , bajo el supuesto de que  $H_0$  es verdadera y que la muestra a su vez tiene una distribución cercana a la normal. Recordemos que  $\bar{x}_N = 527,9$  [ms] y  $s_N = 48$  [ms] en  $n = 1600$  observaciones. Esta distribución se vería como muestra la figura 4.3, creada mediante el script 4.3.

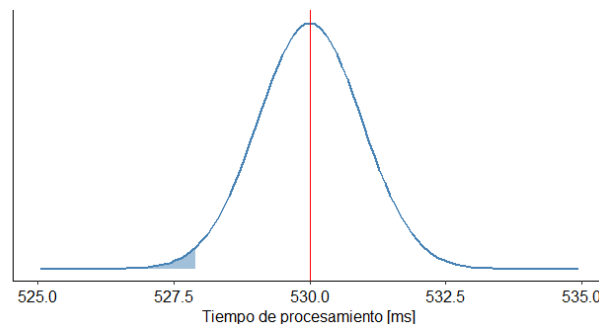


Figura 4.3: probabilidad de encontrar una media igual o menor que  $\bar{x} = 527,9$  [ms] en la distribución muestral con  $\mu_{\bar{x}} = 530$  y  $\sigma_{\bar{x}} = 1,2$ .

En este punto, resulta importante hacer una aclaración en relación al valor p. El área bajo la sección de la curva con valores menores o iguales a un estimador se calcula usando para ello el **valor z**, definido en la ecuación 4.5, como **estadístico de prueba**.

$$z = \frac{\text{estimador puntual} - \text{valor nulo}}{SE_{\text{estimador puntual}}} = \frac{\hat{\theta} - \theta_0}{SE_{\hat{\theta}}} \quad (4.5)$$

Un **estadístico de prueba** es un estadístico de resumen que resulta especialmente útil para evaluar hipótesis o calcular el valor p. El valor z se usa cuando el estimador puntual se acerca a la normalidad, aunque existen otros estadísticos de prueba adecuados para otros escenarios.

Script 4.3: cálculo del valor p para una prueba de una cola.

```

1 library(ggpubr)
2
3 # Generar una muestra donde la media cumpla con la hipótesis nula.
4 set.seed(872)
5
6 media_poblacion_antiguo <- 530
7 media_muestra_nuevo <- 527.9
8 desv_est <- 48
9 n <- 1600
10 error_est <- desv_est / sqrt(n)
11
12 x <- seq(media_poblacion_antiguo - 5.2 * error_est,
13         media_poblacion_antiguo + 5.2 * error_est,
14         0.01)
15
16 y <- dnorm(x, mean = media_poblacion_antiguo, sd = error_est)
17
18 datos <- data.frame(x, y)
19
20 # Graficar la muestra.
21 g <- ggplot(data = datos, aes(x))
22
23 g <- g + stat_function(fun = dnorm,
24                       args = list(mean = media_poblacion_antiguo,
25                                   sd = error_est),
26                       colour = "steelblue", size = 1)
27
28 g <- g + ylab("")
29 g <- g + scale_y_continuous(breaks = NULL)
30 g <- g + scale_x_continuous(name = "Tiempo de procesamiento [ms]")
31 g <- g + theme_pubr()
32
33 # Colorear el área igual o menor que la media observada.
34 g <- g + geom_area(data = subset(datos,
35                                 x < media_muestra_nuevo),
36                   aes(y = y),
37                   colour = "steelblue",
38                   fill = "steelblue",
39                   alpha = 0.5)
40
41 # Agregar una línea vertical para el valor nulo.
42 g <- g + geom_vline(aes(xintercept = media_poblacion_antiguo),
43                    color = "red", linetype = 1)
44
45 print(g)
46
47 # Calcular el valor Z para la muestra.
48 Z <- (media_muestra_nuevo - media_poblacion_antiguo) / error_est
49

```

```

50 # Calcular el valor p.
51 p_1 <- pnorm(Z, lower.tail = TRUE)
52
53 cat("Valor p: ", p_1, "\n")
54
55 # También se puede calcular el valor p directamente a partir de la
56 # distribución muestral definida por el valor nulo y el error
57 # estándar.
58 p_2 <- pnorm(media_muestra_nuevo, mean = media_poblacion_antiguo,
59               sd = est_err)
60
61 cat("Valor p: ", p_2)

```

El valor  $p$ , en este caso  $p = 0,040$ , corresponde al área coloreada en la figura 4.3, y se calcula en la línea 51 del script 4.3. Esto nos indica, en este caso, que si  $H_0$  fuera verdadera y el nuevo sistema tarda en promedio lo mismo que el antiguo en procesar las transacciones, la probabilidad de encontrar una media de a lo más 527,9 [ms] para una muestra de 1.600 transacciones es de 4 %, lo que sería bastante poco frecuente.

Cuanto menor sea el valor  $p$ , más fuerte será la evidencia en favor de  $H_A$  por sobre  $H_0$ . Y aquí la ventaja de usar este método para decidir: el valor  $p$  se puede **comparar directamente** con el nivel de significación  $\alpha$ , y si  $p$  es menor que el nivel de significación se considera evidencia suficiente para rechazar la hipótesis nula en favor de la hipótesis alternativa. En este ejemplo,  $p = 0,040 < \alpha = 0,05$ , por lo que se rechaza  $H_0$  en favor de  $H_A$ . Pero como se dijo cuando usamos intervalos de confianza, el valor  $p$  está cerca del valor  $\alpha$  y convendría ser menos tajante en la decisión y evaluar la posibilidad de ampliar la muestra para conseguir evidencia más definitiva.

Siempre es recomendable formular la conclusión de la prueba de hipótesis en lenguaje llano, para facilitar su comprensión. Así, en este caso concluimos que los datos sugieren que el nuevo sistema tarda menos que el antiguo en procesar transacciones, pero que es necesario hacer un estudio con más observaciones para tener un diagnóstico más definitivo.

Volvamos nuevamente al escenario de la prueba de hipótesis bilateral para el ejemplo, manteniendo el nivel de significación  $\alpha = 0,05$ . Puesto que en este caso nos interesa la diferencia en ambas direcciones, ya que la evidencia en ambas direcciones es favorable para  $H_A$ , debemos considerar el área bajo las dos colas de la curva normal, a diferencia del caso de la prueba de hipótesis unilateral en que solo se consideramos la cola correspondiente a la dirección de interés de la diferencia. Dado que el modelo normal es simétrico, el área bajo ambas colas es la misma (figura 4.4, script 4.4). El valor  $p$ , entonces, ahora es igual a dos veces el área de la cola inferior, es decir,  $p = 0,080$ . Puesto que  $p > \alpha$ , se falla en rechazar  $H_0$ . Es decir, no hay evidencia suficiente para concluir que existe una diferencia entre los tiempos promedio requeridos por ambos sistemas para procesar transacciones.

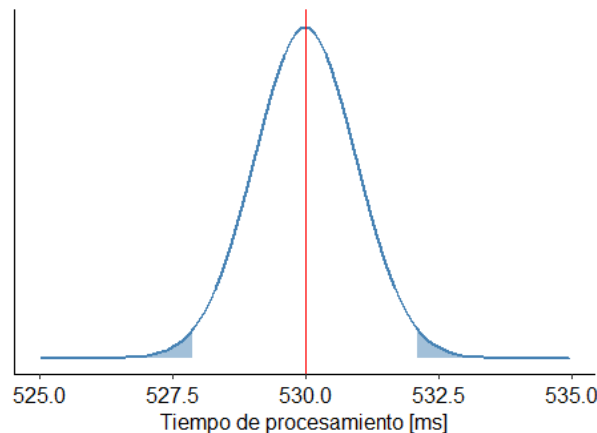


Figura 4.4: cuando la prueba de hipótesis es bilateral, se deben colorear ambas colas.

Script 4.4: cálculo del valor p para una prueba de dos colas.

```
1 library(ggpubr)
2
3 # Generar una muestra donde la media cumpla con la hipótesis nula.
4 set.seed(208)
5
6 media_poblacion_antiguo <- 530
7 media_muestra_nuevo <- 527.9
8 desv_est <- 48
9 n <- 1600
10 error_est <- desv_est / sqrt(n)
11
12 x <- seq(media_poblacion_antiguo - 5.2 * error_est,
13         media_poblacion_antiguo + 5.2 * error_est,
14         0.01)
15
16 y <- dnorm(x,
17           mean = media_poblacion_antiguo,
18           sd = error_est)
19
20 dataframe <- data.frame(x, y)
21
22 # Graficar la muestra.
23 g <- ggplot(data = dataframe, aes(x))
24
25 g <- g + stat_function(fun = dnorm,
26                       args = list(mean = media_poblacion_antiguo,
27                                   sd = error_est),
28                       colour = "steelblue", size = 1)
29
30 g <- g + ylab("")
31 g <- g + scale_y_continuous(breaks = NULL)
32 g <- g + scale_x_continuous(name = "Tiempo de procesamiento [ms]")
33 g <- g + theme_pubr()
34
35 # Colorear el área igual o menor que la media observada.
36 g <- g + geom_area(data = subset(dataframe,
37                                 x < media_muestra_nuevo),
38                   aes(y = y),
39                   colour = "steelblue",
40                   fill = "steelblue",
41                   alpha = 0.5)
42
43 # Calcular el área bajo la cola inferior.
44 area_inferior <- pnorm(media_muestra_nuevo,
45                       mean = media_poblacion_antiguo,
46                       sd = desv_est)
47
48
49 # Colorear igual área en la cola restante.
50 corte_x <- qnorm(1 - area_inferior,
51               mean = media_poblacion_antiguo,
52               sd = desv_est)
53
54 g <- g + geom_area(data = subset(dataframe,
55                                 x > corte_x),
56                   aes(y = y),
57                   colour = "steelblue",
58                   fill = "steelblue",
```

```

59         alpha = 0.5)
60
61 # Agregar una línea vertical para el valor nulo.
62 g <- g + geom_vline(aes(xintercept = media_poblacion_antiguo),
63                      color = "red", linetype = 1)
64
65 print(g)
66
67 # Calcular el valor Z para la muestra.
68 Z <- (media_muestra_nuevo - media_poblacion_antiguo) / error_est
69
70 # Calcular el valor p (recordando ahora que la hipótesis es bilateral).
71 p <- 2 * pnorm(Z, lower.tail = TRUE)
72
73 cat("Valor p: ", p)

```

Un punto importante que debemos tener en cuenta es que **las pruebas unilaterales** se usan cuando se desea verificar un incremento o un decremento, pero no ambas. No obstante, esta decisión debe tomarse siempre **antes de examinar los datos**, pues de lo contrario se duplica la probabilidad de cometer errores de tipo I y se está cayendo en **prácticas poco éticas**.

#### 4.5.2 El efecto del nivel de significación

Hemos visto que el nivel de significación ( $\alpha$ ) representa la proporción de veces en que se cometería un error de tipo I (es decir, rechazar  $H_0$  en favor de  $H_A$ , cuando  $H_0$  es en realidad verdadera). Si resulta costoso o peligroso cometer un error de este tipo, debemos requerir evidencia más fuerte para rechazar la hipótesis nula (es decir, reducir la probabilidad de que esto ocurra), lo que podemos lograr usando un valor más pequeño para el nivel de significación, por ejemplo,  $\alpha = 0,01$ . Sin embargo, esto necesariamente **aumentará** la probabilidad de cometer un error de tipo II.

Si, por el contrario, el costo o el peligro de cometer un error de tipo II (no rechazar  $H_0$  cuando en realidad  $H_A$  es verdadera) es mayor, debemos escoger un nivel de significación más elevado (por ejemplo,  $\alpha = 0,10$ ).

Así, **el nivel de significación seleccionado para una prueba siempre debe reflejar las consecuencias de cometer errores de tipo I o de tipo II**.

## 4.6 INFERENCIA PARA OTROS ESTIMADORES

Hasta ahora, solo hemos considerado la media como estimador para la inferencia. No obstante, muchos de los conceptos que hemos visto en este capítulo pueden aplicarse, con algunas ligeras modificaciones, usando otros estimadores.

#### 4.6.1 Estimadores puntuales con distribución cercana a la normal

En realidad existen múltiples estimadores puntuales, además de la media, cuya distribución muestral es cercana a la normal si las muestras son lo suficientemente grandes, tales como las proporciones y la diferencia de medias. Si bien veremos con detalle la prueba de hipótesis con estos estimadores puntuales en capítulos posteriores, es importante contar con algunas orientaciones generales.

Un supuesto importante que debemos tener en cuenta es que el estimador puntual  $\hat{\theta}$  debe ser **insesgado**. Esto significa que la distribución muestral de  $\hat{\theta}$  tiene su centro en el valor del parámetro  $\theta$  que estima. En otras palabras, un estimador insesgado (como la media) tiende a proveer una estimación cercana al parámetro real.

En términos generales, el intervalo de confianza para un estimador puntual insesgado cuya distribución es cercana a la normal (como la media, las proporciones o la diferencia de medias) está dado por la ecuación 4.6, donde  $z^*$  se escoge de manera tal que se condiga con el nivel de confianza seleccionado y y la lateralidad de la hipótesis alternativa. Como se dijo anteriormente, el valor  $z^* \cdot SE_{\hat{\theta}}$  se denomina “margen de error”. Debemos recordar que la ecuación 4.2 corresponde al error estándar de la media, pero los errores estándar para otros estimadores puntuales se estiman de manera diferente a partir de los datos.

$$\hat{\theta} \pm z^* \cdot SE_{\hat{\theta}} \quad (4.6)$$

El método de prueba de hipótesis usando valores p puede generalizarse para otros estimadores puntuales con distribución cercana a la normal. Para ello, Diez y col. (2017, p. 199) señalan que se debemos considerar los siguientes pasos:

Prueba de hipótesis usando el modelo normal:

1. Formular las hipótesis nula ( $H_0$ ) y alternativa ( $H_A$ ) en lenguaje llano y luego en notación matemática.
2. Identificar un estimador puntual (estadístico) adecuado e insesgado para el parámetro de interés.
3. Verificar las condiciones para garantizar que la estimación del error estándar sea razonable y que la distribución muestral del estimador puntual siga aproximadamente una distribución normal.
4. Calcular el error estándar. Luego, graficar la distribución muestral del estadístico bajo el supuesto de que  $H_0$  es verdadera y sombrear las áreas que representan el valor p.
5. Usando el gráfico y el modelo normal, calcular el valor p para evaluar las hipótesis y escribir la conclusión en lenguaje llano.

#### 4.6.2 Estimadores con otras distribuciones

Existen métodos de construcción de intervalos de confianza y prueba de hipótesis adecuados para aquellos casos en que el estimador puntual o el estadístico de prueba no son cercanos a la normal (por ejemplo, si la muestra es pequeña, se tiene una mala estimación del error estándar o el estimador puntual tiene una distribución distinta a la normal). No obstante, la selección de métodos alternativos debe hacerse siempre teniendo en cuenta la distribución muestral del estimador puntual o del estadístico de prueba.

Una consideración importante es que **siempre debemos verificar el cumplimiento de las condiciones requeridas por una herramienta estadística**, pues de lo contrario las conclusiones pueden ser erradas y carecerán de validez.

## 4.7 EJERCICIOS PROPUESTOS

1. ¿Es correcto afirmar que, si se lanza un dado una y otra vez, la media móvil simple del número de puntos que aparecen en la cara superior crece monótonamente? Justifica tu respuesta.
2. ¿Es correcto afirmar que, si se lanza un dado una y otra vez, la proporción de veces que aparece un número impar de puntos (1, 3 o 5) en la cara superior es siempre 0,5? Justifica su respuesta.
3. Si se calcula la media de diez muestras distintas extraídas de la misma población, ¿se espera ver el mismo valor cada vez? ¿Cómo se llama a este fenómeno?
4. Completa las siguientes oraciones:
  - a) Una estimación \_\_\_\_\_ es un \_\_\_\_\_ calculado con datos de una muestra como aproximación del valor desconocido de un \_\_\_\_\_ de la población en estudio.
  - b)  $\bar{X}$  o  $\bar{x}$  se usan para denotar la \_\_\_\_\_, que es una estimación puntual de  $\mu$ , la \_\_\_\_\_.
5. Se sabe que una prueba para medir el coeficiente intelectual de jóvenes de 18 años produce puntuaciones que siguen una distribución  $\mathcal{N}(\mu = 100, \sigma^2 = 100)$ .
  - a) Dibuja el histograma de la distribución muestral de medias para muestras de tamaño 25 de esta población.
  - b) Una de las muestras anteriores presentó  $\bar{x} = 95$  y  $s = 15$ . Determina el intervalo con 95 % de confianza para este caso.
  - c) Con otra de las muestras se pudo determinar que su intervalo con 99 % confianza era  $[90, 26; 105, 74]$ . ¿Qué significa esto?
  - d) El intervalo anterior, ¿es más grande o más pequeño que uno con 90 % de confianza?
6. Una empresa de tecnología quiere promocionar un software especializado para almacenar y recuperar imágenes médicas digitales. Con esta idea, está financiando un estudio para determinar el tiempo (en segundos) que necesita un grupo de médicos para recuperar imágenes desde sus propios registros en sus portátiles personales y desde la base de datos central con el software ofrecido y una conexión a la Web.
  - a) Enuncia las hipótesis nula y alternativa (en castellano común).
  - b) Identifica la variable aleatoria que se va a estudiar, el parámetro de interés y el correspondiente estadístico.
  - c) Enuncia, más formalmente, las hipótesis nula y alternativa para este caso.
  - d) Supón que el intervalo con 95 % confianza para el tiempo de recuperación promedio de una imagen digital desde la base de datos central resultó ser  $[24; 36]$  [s]. ¿Qué decisión tomarías ante la hipótesis nula: la media del tiempo de recuperación de una imagen digital con el nuevo software es de 25 segundos? En este caso, ¿cuál podría ser la hipótesis alternativa?
  - e) Para el intervalo de confianza anterior, ¿cuál sería un error de tipo I?
  - f) Conociendo el intervalo de confianza anterior, ¿es posible cometer un error de tipo II? Explica.
7. Si una hipótesis nula es falsa, aumentar el nivel de significación para un tamaño de muestra dado, ¿reduce la probabilidad de rechazarla?
8. ¿Qué significa que un estadístico tenga un valor p de 0,025?
9. Si una hipótesis nula es rechazada a un nivel de significación de 0,01, ¿será rechazada a un nivel de significación 0,05? Explica.
10. Si una hipótesis nula es rechazada por una prueba unilateral (una cola), ¿será también rechazada por una prueba bilateral (dos colas)? Explica.
11. Acabas de leer un artículo que hace la siguiente aseveración: “a 95 % confidence interval for mean reaction time is from 0.25 to 0.29 seconds. Thus, about 95 % of individuals will have reaction times in this interval.” Comenta.
12. Da el ejemplo de un estudio en que es más dañino cometer un error tipo II que un error tipo I.
13. Lista las condiciones que deben verificarse para asegurar que el TLC (teorema del límite central) está rigiendo y es posible hacer una prueba de hipótesis o calcular un intervalo de confianza.
14. Si para un estudio de una determinada variable aleatoria numérica es igualmente dañino cometer errores de tipo I como errores tipo II:
  - a) Dibuja la distribución de una muestra de tamaño 16 (un diagrama de caja, por ejemplo) para la que el contraste de hipótesis con nivel de significación 0,05 sea confiable.



- b)* Dibuja la distribución de una muestra de tamaño 30 en que se requiera de un nivel de significación más exigente ( $\alpha < 0,05$ ) para hacer el contraste de hipótesis más confiable.
  - c)* Dibuja la distribución de una muestra en que es mejor no confiar en el contraste de hipótesis con métodos estudiados hasta ahora.
- 15. Si un estudio sobre el tiempo promedio de búsqueda y recuperación de imágenes médicas con dos tecnologías distintas reporta: “existe una diferencia significativa ( $p < 0,02$ ) entre el tiempo invertido con la tecnología A ( $33 \pm 4[s]$ ) que con la tecnología B ( $30 \pm 6[s]$ )”, ¿significa que se debe adoptar la tecnología B? ¿Por qué?
- 16. Explica por qué se incrementa la probabilidad de cometer errores tipo I al cambiar de una prueba de hipótesis bilateral a otra unilateral.



## REFERENCIAS

- Bache, S. (2014). *Introducing magrittr*. Consultado el 7 de abril de 2021, desde <https://cran.r-project.org/web/packages/magrittr/vignettes/magrittr.html>
- Carchedi, N., De Mesmaeker, D. & Vannoorenberghe, L. (s.f.). RDocumentation. Consultado el 2 de abril de 2021, desde <https://www.rdocumentation.org/>
- Cross, J. (2017). *Discrete Random Variables*. Consultado el 9 de abril de 2021, desde <https://rpubs.com/jcross/discreteRV>
- Dagnino, J. (2014). Tipos de datos y escalas de medida. *Revista Chilena de Anestesia*, 42(2), 109-111.
- Devore, J. L. (2008). *Probabilidad y Estadística para Ingeniería y Ciencias* (7.ª ed.). CENAGE Learning.
- Diez, D., Barr, C. D. & Çetinkaya-Rundel, M. (2017). *OpenIntro Statistics* (3.ª ed.). <https://www.openintro.org/book/os/>.
- Field, A., Miles, J. & Field, Z. (2012). *Discovering statistics using R*. SAGE Publications Ltd.
- Freedman, D. A. (2009). *Modelización*. Cambridge University Press.
- Freund, R. J. & Wilson, W. J. (2003). *Statistical Methods* (2.ª ed.). Academic Press.
- Joly, F. (1988). *La Cartografía*. Oikos-Tau, S.A. Ediciones.
- Kaplan, D. (2009). *Statistical Modeling: A Fresh Approach*. Consultado el 8 de marzo de 2019, desde [http://works.bepress.com/daniel\\_kaplan/38](http://works.bepress.com/daniel_kaplan/38)
- McCullagh, P. (2002). What Is a Statistical Model? *The Annals of Statistics*, 30(5), 1225-1267.
- Mendez Ramírez, I. (1998). Empirismo, método científico y estadística. *Revista de Geografía Agrícola (Mexico)*.
- Mendez Ramírez, I. (2012). Método Científico: aspectos epistemológicos y metodológicos para el uso de la Estadística. *SaberEs*, 4.
- Müller, K. (2021). *dplyr*. Consultado el 10 de septiembre de 2021, desde <https://dplyr.tidyverse.org/>
- Real Academia Española. (2014). *Diccionario de la lengua española* (23.ª ed.). Consultado el 30 de marzo de 2021, desde <https://dle.rae.es>
- Ríos, S. (1995). *Modelización*. Alianza Ediciones.
- RStudio. (2021). RStudio. Consultado el 2 de abril de 2021, desde <https://rstudio.com/products/rstudio/>
- SAS Institute Inc. (2008). SAS/STAT® 9.2 User's Guide.
- STDHA. (s.f.). Descriptive Statistics and Graphics - Easy Guides - Wiki - STHDA. Consultado el 31 de marzo de 2021, desde <http://www.sthda.com/english/wiki/descriptive-statistics-and-graphics>
- The R Foundation. (s.f.). Documentation. Consultado el 2 de abril de 2021, desde <https://www.r-project.org/other-docs.html>
- Wickham, H. (2021). *tidyr*. Consultado el 10 de septiembre de 2021, desde <https://r4ds.had.co.nz/index.html>
- Wickham, H. & Grolemond, G. (2017). *R for Data Science*. <https://r4ds.had.co.nz/index.html>
- Willems, K. (2017). *Formulas in R Tutorial*. Consultado el 11 de septiembre de 2021, desde <https://dplyr.tidyverse.org/>