

"ESTADÍSTICA PARA LAS CIENCIAS SOCIALES", FERRIS S. RITCHEY

CAPÍTULO

5

MEDICIÓN DE LA DISPERSIÓN O VARIACIÓN EN UNA DISTRIBUCIÓN DE PUNTUACIONES

Introducción	126
El rango	127
Limitaciones del rango: situaciones en las que reportarlo solo puede conducir a errores	129
La desviación estándar	129
Pensamiento proporcional y lineal sobre la desviación estándar	130
Método abreviado para calcular la desviación estándar	135
Limitaciones de la desviación estándar	136
La desviación estándar como parte integral de la estadística inferencial	137
¿Por qué se llama desviación "estándar"?	137
Puntuaciones estandarizadas (puntuaciones Z)	138
La distribución normal	140
Uso del rango para estimar la desviación estándar	143
Una ilustración completa del cálculo de los estadísticos de dispersión	144

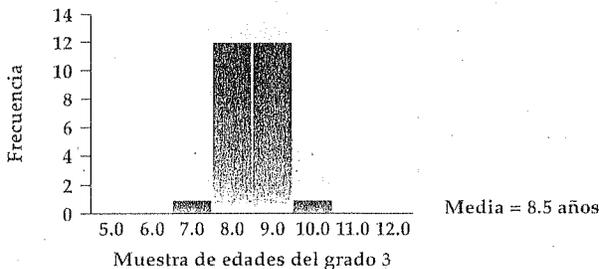
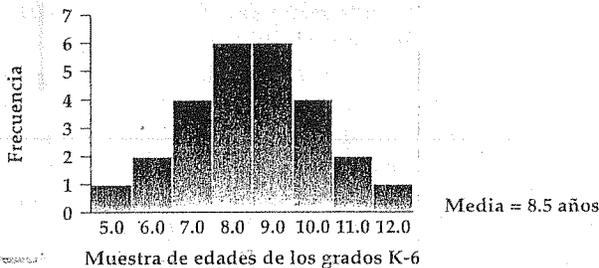
Uso de un formato desglosado para calcular la desviación estándar	144
Uso de un formato de distribución de frecuencias para calcular la desviación estándar	145
Presentación tabular de resultados	146
Insensatez y falacias estadísticas: ¿qué indica cuando la desviación estándar es más grande que la media?	147

Introducción

Reportar un estadístico de tendencia central por sí mismo no es suficiente para comunicar la forma de una distribución de puntuaciones. Dos muestras con las mismas medias pueden tener formas sumamente diferentes. La figura 5-1 presenta dos distribuciones de edades: para una muestra de alumnos de escuela primaria (desde jardín de niños hasta sexto grado, o K-6) y para una clase de tercer grado de una segunda escuela. La edad media de los alumnos en ambas escuelas es de 8.5 años. En la escuela K-6, sin embargo, los niños tienen entre 5 y 12 años. En la clase de tercer grado de la otra escuela ninguno de los alumnos es menor de 7 o mayor de 10 años de edad. Aunque estas dos distribuciones de edades tienen la misma tendencia central, sus puntuaciones se dispersan de manera muy diferente, con una mayor dispersión de edades en la escuela K-6.

FIGURA 5-1

Comparación de la dispersión de edades de alumnos en dos muestras con las mismas medias



El tema de este capítulo es la **dispersión**, es decir, *cómo se extienden las puntuaciones de una variable de intervalo/razón de la menor a la mayor y la forma de la distribución entre éstas*. (Como una ayuda para la memoria, recuerde que Johnny Appleseed *dispersó* semillas de manzana.) Existe un número infinito de posibles formas de distribución para una variable con una media dada. Todas las puntuaciones podrían agruparse alrededor de la media con la clara forma de una curva de campana; pero la curva podría ser de diferentes tamaños, dependiendo del tamaño de la muestra. O las puntuaciones podrían estar ligera o grandemente sesgadas hacia un lado. Además, de esto una sola variable puede tener diferentes dispersiones de una población a otra. Por ejemplo, el ingreso familiar anual para residentes en Estados Unidos varía desde cero hasta decenas de millones de dólares; mientras el ingreso familiar de los pobres que viven en proyectos de asistencia social varía desde cero hasta unos cuantos miles de dólares.

Dispersión

Cómo se extienden las puntuaciones de una variable de intervalo/razón de la menor a la mayor y la forma de la distribución entre éstas.

Los estadísticos de **dispersión** describen *cómo se extienden las puntuaciones de una variable de intervalo/razón a través de su distribución*. Los estadísticos de dispersión permiten descripciones precisas de la frecuencia de casos en cualquier punto de una distribución. Por ejemplo, si el gobierno federal decide aumentar los impuestos para los "ricos", empleando estadísticos de dispersión podemos identificar el nivel de ingresos del 5 por ciento más rico de todas las familias del país. De manera similar, si un programa de asistencia social se planea para cubrir sólo 10 000 familias de la ciudad, podemos establecer qué nivel de ingreso familiar califica para recibir la asistencia. Estudiar la dispersión es como pasear hacia atrás y hacia adelante a lo largo del eje X de un histograma, y observar dónde se concentran los casos. ¿La mayoría de los casos cae alrededor de la media o están cargados hacia algún lado? ¿Cuántos casos caen entre cualesquiera dos puntos? ¿Qué valor de la variable corta el 10 por ciento superior de casos? Los dos estadísticos de dispersión más usados se analizan más adelante: el rango y la desviación estándar.

Estadísticos de dispersión

Estadísticos que describen cómo se extienden las puntuaciones de una variable de intervalo/razón a través de su distribución.

El rango

El rango es una expresión de *cómo las puntuaciones de una variable de intervalo/razón se distribuyen de la menor a la mayor* —la distancia entre las puntuaciones mínima y

máxima en una muestra—. Se calcula como la diferencia entre las puntuaciones máxima y mínima, más el valor de la unidad de redondeo. El valor de la unidad de redondeo (por ejemplo, 1 si las puntuaciones se redondean al número entero más cercano, 0.1 si las puntuaciones se redondean a la décima más cercana y así sucesivamente) se suma para considerar el límite real inferior de la puntuación más baja y el límite real superior de la puntuación más alta.

Cálculo del rango de una variable X de intervalo/razón

1. Ordene las puntuaciones en la distribución de menor a mayor.
2. Identifique las puntuaciones mínima y máxima.
3. Identifique el valor de la unidad de redondeo (véase el apéndice A como repaso).
4. Calcule el rango:

$$\text{Rango} = (\text{puntuación máxima} - \text{puntuación mínima}) + \text{valor de la unidad de redondeo}$$

El rango

Expresión de cómo las puntuaciones de una variable de intervalo/razón se distribuyen de menor a mayor.

Calculemos el rango para un problema de ejemplo. Suponga que X = edad (redondeada al año más cercano) y tenemos la siguiente distribución de puntuaciones:

21, 23, 43, 26, 20, 21, 25

Empiece por ordenar las puntuaciones:

20, 21, 21, 23, 25, 26, 43

Identifique las puntuaciones mínima y máxima de 20 y 43, respectivamente, e identifique que la unidad de redondeo es 1.

Calcule el rango:

$$\begin{aligned} \text{Rango} &= (\text{puntuación máxima} - \text{puntuación mínima}) + \text{valor de la unidad de redondeo} \\ &= (43 - 20) + 1 = 24 \text{ años} \end{aligned}$$

Como resultado del redondeo, el individuo que registró 20 podría tener 19.5 años; y el de 43 años, podría tener 43.5 años. El rango de 24 años es la distancia entre estos límites reales menor y mayor de las puntuaciones; es decir, 43.5 años - 19.5 años = 24 años.

A menudo resulta más informativo reportar las puntuaciones mínima y máxima por sí mismas, señalando que estas edades varían desde 20 hasta 43. De esta manera indirectamente mostramos que en la muestra no hay menores de 20 años ni mayores de 43 años de edad.

Limitaciones del rango: situaciones en las que reportarlo solo puede conducir a errores

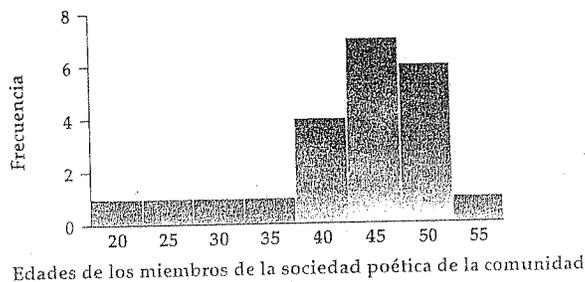
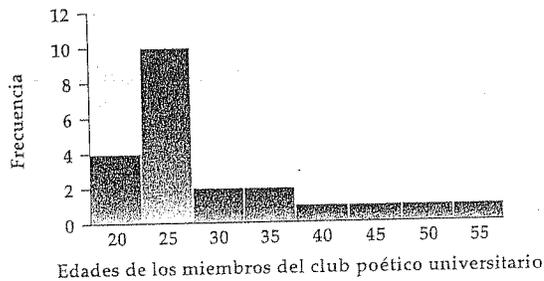
Puesto que el rango utiliza las puntuaciones más extremas de una distribución, un valor extremo inflará enormemente su cálculo. Esto sucedió para las siete edades indicadas arriba. Los 43 años hicieron que el rango pareciera estar extendido por encima de los 24 años. Reportar esto daría la impresión de que la muestra tiene un número considerable de sujetos de 30 y 40 años. Un reporte más exacto estipularía que con la excepción del estudiante de 43 años, las edades tenían un rango de 7 años ($26 - 20 + 1 = 7$). Omitir el valor extremo e indicarlo como una excepción es una forma razonable de ajustar esta limitación del rango.

El rango también está limitado por su estrecho alcance informativo. No nos dice nada sobre la forma de la distribución entre las puntuaciones extremas. Por ejemplo, las dos distribuciones ilustradas en la figura 5-2 tienen el mismo rango, sugiriendo formas similares; pero, de hecho, sus formas son radicalmente diferentes. Por último, hay poco que uno pueda hacer matemáticamente con el rango. En suma, el rango tiene utilidad limitada, sobre todo cuando se informa solo.

La desviación estándar

La desviación estándar es otra medida sumaria de la dispersión o la variación de las puntuaciones en una distribución. Este estadístico de dispersión es muy dife-

FIGURA 5-2
Comparación de dos distribuciones con formas diferentes que tienen el mismo rango



rente del rango. Al enfocarse en los extremos de la distribución, el rango se aproxima a la dispersión desde los "exteriores" o extremos de la distribución. Observar el rango es como mirar un juego de básquetbol desde lo alto de las tribunas; la cancha parece encajonada por las canastas en cada extremo. En contraste, la *desviación estándar describe cómo las puntuaciones una variable de intervalo/razón u ordinal de tipo intervalo se extienden a lo largo de la distribución, en relación con la puntuación media*. La media es un estadístico de tendencia central y como tal proporciona un punto de enfoque que se centra "dentro" de la distribución. Observar la dispersión a partir de la media con su desviación estándar es como mirar desde el centro de la cancha; el centro de atención está en la distancia desde el centro de la cancha hasta otros puntos en cualquier dirección. Como la media, la desviación estándar es muy apropiada con variables de intervalo/razón y variables ordinales de tipo intervalo.

Desviación estándar

Describe cómo las puntuaciones de una variable de intervalo/razón u ordinal de tipo intervalo se extienden a lo largo de la distribución, en relación con la puntuación media.

Pensamiento proporcional y lineal sobre la desviación estándar

La desviación estándar se calcula determinando qué tan lejos está cada puntuación de la media —qué tan lejos se *desvía* de la media—. En este sentido, la desviación estándar es un derivado (o producto) de la media, y las dos medidas siempre se informan juntas. De hecho, la frase "la media y la desviación estándar" es una de las más empleadas por los estadísticos. La desviación estándar —como una medida sumaria de todas las puntuaciones en una distribución— nos dice con qué amplitud se agrupan las puntuaciones alrededor de la media. Como brevemente lo analizaremos, la desviación estándar también es útil en conjunción con la curva normal.

Antes de calcular una desviación estándar, analicemos su fórmula. La siguiente fórmula sirve para calcular directamente la desviación estándar:

Método directo para calcular la desviación estándar

$$s_x = \sqrt{\frac{\sum(X - \bar{X})^2}{n - 1}}$$

Donde

s_x = desviación estándar para la variable X de intervalo/razón

\bar{X} = media de X

n = tamaño de la muestra

Vale la pena aproximarnos paso a paso al cálculo de la desviación estándar. Esto elimina el misterio de la fórmula (con su Σ , cuadrada, y símbolos de raíz cuadrada) y nos ayuda a apreciar que la desviación estándar es parte esencial de la curva normal.

Identifique las especificaciones. Empezamos identificando la información dada. Especificación: X = una variable de intervalo/razón (u ordinal de tipo intervalo), n = tamaño de la muestra, y la distribución de puntuaciones en bruto para X .

Calcule la media. Calculamos la media porque la desviación estándar está diseñada para medir la dispersión alrededor la media.

$$\bar{X} = \frac{\Sigma X}{n}$$

Calcule las puntuaciones de la desviación: pensamiento lineal. Luego determinamos qué tan lejos de la media cae la puntuación de cada sujeto. La diferencia entre una puntuación y su media se llama **puntuación de desviación**, es decir, *cuánto difiere o se "desvía" de la media una puntuación individual*:

$$X - \bar{X} = \text{puntuación de desviación para un valor de } X$$

Piense en una puntuación de desviación como una medida de distancia en el eje X . ¿Qué nos dice la puntuación de desviación? Suponga que X es la variable peso, y el peso medio para una muestra de jugadoras de voleibol de la Universidad de Elmstown es de 138 libras. La jugadora estrella, Sandra "Clavadora de alma" Carson, pesa 173 libras; ésa es su puntuación en bruto o "puntuación X ". Su puntuación de desviación es más 35 libras:

$$\text{Puntuación de desviación} = X - \bar{X} = 173 - 138 = 35 \text{ libras}$$

La puntuación de desviación nos dice dos cuestiones sobre una puntuación en la distribución: 1) la magnitud de la distancia desde la puntuación X hacia la media, y 2) la dirección de la puntuación X : ya sea que esté abajo o arriba de la media. Cuando una puntuación X es mayor que la media, la puntuación de desviación resultará un valor positivo, como el de Sandra, lo cual significa que la puntuación X queda a la derecha en una curva de distribución. Cuando una puntuación X es menor que la media, la puntuación de desviación resultará negativa, lo que significa que la puntuación X queda a la izquierda de la media. La puntuación de desviación de Sandra de +35 libras nos indica que ella está 35 libras *encima* del peso medio del equipo.

Puntuación de desviación

Cuánto difiere o "se desvía" de la media una puntuación individual.

La puntuación de desviación es el cálculo matemático central al determinar la desviación estándar. Como una medida sumaria para la muestra entera, la desviación estándar es una suma y promedio del cuadrado de estas puntuaciones de desviación, como en los siguientes pasos.

Suma las puntuaciones de desviación. El siguiente paso para calcular la desviación estándar consiste en sumar las puntuaciones de desviación. Tal suma siempre será igual a cero (dentro del error de redondeo):

$$\Sigma(X - \bar{X}) = 0 = \text{suma de las puntuaciones de desviación}$$

La suma de puntuaciones de desviación constituye una verificación respecto de la exactitud de los cálculos, porque la suma de puntuaciones de desviación *siempre* será igual a cero (con cierto error de redondeo). En el capítulo 4 estudiamos cómo la media es un punto de balance en la distribución. Lo que la media hace es balancear las desviaciones, para que se cancelen entre sí y resulten en una suma de puntuaciones de desviación igual a cero. De hecho, otra definición matemática de la *media es aquel punto en una distribución donde las puntuaciones de desviación suman cero.*

Eleve al cuadrado las puntuaciones de desviación y sume los cuadrados. La dispersión de una variable a menudo se compara para dos o más muestras. El hecho de sumar las puntuaciones de desviación no detectará una diferencia en la dispersión entre dos muestras, porque la suma para ambas será cero. Esto potencialmente nos deja en un callejón sin salida. Si las puntuaciones de una muestra se dispersan ampliamente y las de la otra lo hacen de manera estrecha, ¿qué beneficio implica informar que ambas tienen una suma de puntuaciones de desviación de cero? ¡Ninguno! Por consiguiente, al comparar dos muestras, debemos encontrar una manera de sumar las puntuaciones de desviación para que la suma sea más grande para una muestra con una dispersión mayor. La solución más útil consiste en elevar al cuadrado cada puntuación de desviación y después sumar los cuadrados. Al elevar al cuadrado se eliminan los signos negativos en las puntuaciones de desviación. La *suma de las puntuaciones de desviación al cuadrado es la variación* (a menudo se denomina *suma de cuadrados*), un estadístico que resume las desviaciones para la muestra entera:

$$\Sigma(X - \bar{X})^2 = \text{la variación (o "suma de cuadrados")}$$

La variación o la suma de cuadrados

Es la suma de las puntuaciones de desviación al cuadrado; un estadístico que resume las desviaciones para la muestra entera.

Divida la suma de cuadrados entre $n - 1$ para ajustar el tamaño y el error de la muestra: pensamiento proporcional. La suma de cuadrados, o variación, constituye una buena medida de la dispersión de una distribución; pero este estadístico presenta dos problemas. Primero, suponga que deseamos comparar las distribuciones de dos muestras de tamaños diferentes. Por ejemplo, podemos comparar las dis-

tribuciones de los promedios para muestras de estudiantes de la universidad local ($n = 88$) y de la universidad estatal ($n = 104$). Cuando sumamos los cuadrados para cada muestra, podría suceder que obtuviéramos una suma más alta para la universidad estatal, simplemente porque sumamos más números —104 casos en lugar de tan sólo 88—. Cada puntuación X suma alguna cantidad al cálculo. En otras palabras, todo lo demás es igual, cuanto más observaciones existan, mayor será la suma de cuadrados. Para realizar una comparación equilibrada de dos muestras de tamaño diferente, entonces, necesitamos ajustar el número de observaciones en cada muestra dividiendo cada uno entre su tamaño muestral (n). Esto nos da la *variación promedio* (la media de la suma de cuadrados) en cada muestra. De esta forma ajustamos la suma de cuadrados en proporción al número de casos en la muestra.

Una segunda consideración respecto al tamaño de la muestra es que incluirá al error de muestreo; cuanto mayor sea la muestra, menor será el error de muestreo. Los estadísticos han determinado que si restamos 1 de n , este pequeño ajuste produce un estadístico de la muestra que estima con mayor precisión el parámetro de la población. Simplemente, restando 1 del tamaño de la muestra, realizamos un ajuste para el error de muestreo. (Considere que con muestras grandes, este ajuste tendría poco efecto en el cálculo; mientras que con muestras pequeñas, tendría un gran efecto.)

En resumen, dividimos la *variación* (suma de cuadrados) entre $n - 1$ para considerar tanto los efectos del tamaño de la muestra en la suma como el error de muestreo. El resultado se llama *varianza*, y su símbolo es s_x^2 :

$$s_x^2 = \frac{\Sigma(X - \bar{X})^2}{n - 1} = \text{varianza de una muestra}$$

La *varianza* es la *variación promedio de las puntuaciones en una distribución*. Para evitar confundir *varianza* y *variación*, note el sonido acentuado en “*varianza*” y note que n está en su denominador. (Finalmente, debemos hacer notar que si la *desviación estándar* se calcula para las puntuaciones de una población entera, el error de muestreo no constituiría un problema. Por consiguiente, no necesitamos restar 1 de n para obtener la *variación* de una población, lo que se simbolizaría como σ_x^2 .)

Varianza

Variación promedio de las puntuaciones en una distribución (es decir, la media de la suma de cuadrados).

Saque la raíz cuadrada de la *varianza* para obtener la *desviación estándar*. Para producir una buena medida de dispersión se requiere un último paso. La *varianza* es absolutamente aceptable para cálculos; pero no se interpreta de manera directa porque las unidades de medida están elevadas al cuadrado. Así, podríamos calcular la *varianza* de peso para el equipo de fútbol de la universidad local, y encontraríamos que es de 1 391.45 libras cuadradas. No obstante, ¿qué es una “libra cuadrada”? Es una libra de veces una libra; pero ¿quién sabe lo que realmente significa, excepto quizás un matemático? Necesitamos una unidad de medida directamente interpretable —libras en lugar de libras al cuadrado. Para “regresar” a libras, saca-

mos la raíz cuadrada de la varianza. (La raíz cuadrada de una unidad de medida al cuadrado es la unidad de medida en sí.) El resultado es la desviación estándar:

$$s_x = \sqrt{\frac{\Sigma(X - \bar{X})^2}{n - 1}} = \sqrt{s_x^2}$$

En el caso del peso del equipo local, la desviación estándar sería 37.30 libras:

$$s_x = \sqrt{\frac{\Sigma(X - \bar{X})^2}{n - 1}} = \sqrt{s_x^2}$$

$$= \sqrt{1391.45} = 37.30 \text{ libras}$$

Los pasos previamente esbozados involucran un cálculo directo de la desviación estándar. Los elementos de la ecuación—las puntuaciones de desviación, la suma de cuadrados o variación y la varianza—son importantes por sí mismos. Estos elementos aparecen por sí mismos en muchas fórmulas estadísticas (véase, por ejemplo, el capítulo 12). Los pasos para calcular directamente la desviación estándar se resumen en la tabla 5-1, que usted encontrará muy útil para los capítulos posteriores.

TABLA 5-1 Comprensión de la desviación mediante su cálculo directo

Pasos para calcular la desviación estándar	Lo que el paso logra
1. Identifique las especificaciones	1. X debe ser una variable de intervalo/razón (u ordinal de tipo intervalo).
2. Calcule la media: $\bar{X} = \frac{\Sigma X}{n}$	2. Puesto a que la desviación estándar está basada en desviaciones de la media.
3. Calcule las puntuaciones de desviación: $X - \bar{X}$	3. Para determinar la distancia de cada puntuación hacia la media.
4. Sume las puntuaciones de desviación $\Sigma(X - \bar{X})$	4. Asegúrese de que $\Sigma(X - \bar{X}) = 0$
5. Eleve al cuadrado las puntuaciones de desviación y súmelas para obtener la variación o suma de cuadrados: $\text{Variación} = \Sigma(X - \bar{X})^2$	5. Las puntuaciones de desviación se elevan al cuadrado para eliminar los signos negativos y obtener una suma diferente de cero.
6. Calcule la varianza: $s_x^2 = \frac{\Sigma(X - \bar{X})^2}{n - 1}$	6. Divida la suma de cuadrados entre $n - 1$ para ajustarla para el tamaño de la muestra y el error de muestreo.
7. Calcule la desviación estándar, s_x : $s_x = \sqrt{\frac{\Sigma(X - \bar{X})^2}{n - 1}} = \sqrt{s_x^2}$	7. Saque la raíz cuadrada de la varianza para obtener unidades de medida directamente interpretables (unidades en lugar de unidades al cuadrado).

89

TABLA 5-2 Formato desglosado para calcular la desviación estándar usando los métodos directo y abreviado: peso de los jugadores de fútbol de la universidad local ($n = 12$)

Especificaciones		Cálculos		
(1) Jugador	(2) X	(3) $X - \bar{X}$	(4) $(X - \bar{X})^2$	(5) X^2
1	165	-73		
2	200	-38	5 329	27 225
3	216	-22	1 444	40 000
4	217	-21	484	46 656
5	226	-12	441	47 089
6	236	-2	144	51 076
7	239	1	4	55 696
8	244	6	1	57 121
9	261	23	36	59 536
10	268	30	529	68 121
11	283	45	900	71 824
12	301	63	2 025	80 089
$n = 12$	$\Sigma X = 2 856$ libras	$\Sigma(X - \bar{X}) = 0$	$\Sigma(X - \bar{X})^2 = 15 306$ libras cuadradas	$\Sigma X^2 = 695 034$ libras cuadradas

Resulta una buena práctica preparar una tabla desglosada para estos cálculos. La tabla 5-2 presenta una tabla desglosada para determinar la desviación estándar de los pesos de 12 de los 98 jugadores en el equipo de fútbol local. (La columna 5 de la tabla se emplea con un método de cálculo abreviado que se describirá brevemente.) Para calcular las puntuaciones de desviación, $X - \bar{X}$, calculamos la media y le restamos cada puntuación para obtener la tercera columna de la tabla desglosada:

$$\bar{X} = \frac{\Sigma X}{n} = \frac{2 856}{12} = 238 \text{ libras}$$

Finalmente, elevamos al cuadrado las puntuaciones de desviación en la columna 3 para obtener la columna 4. La suma en la columna 4 de la tabla 5-2 y el tamaño de la muestra n son todo lo que necesitamos para calcular la desviación estándar:

$$s_x = \sqrt{\frac{\Sigma(X - \bar{X})^2}{n - 1}} = \sqrt{\frac{15 306}{11}} = \sqrt{1 391.45} = 37.30 \text{ libras}$$

Método abreviado para calcular la desviación estándar

Nuestro propósito para realizar el cálculo directo de la desviación estándar consiste en desarrollar un sentido de proporción acerca de cómo mide las desviaciones con respecto a la media. Sin embargo, si ésta fuera una tarea diaria, calcular direc-

tamente la desviación estándar resultaría embarazoso y propenso al error. Requiere calcular la media y, después, efectuar muchas restas también propensas al error, algo que es especialmente fastidioso con números decimales. Por fortuna, existe un método abreviado de cálculo al cual se llega sustituyendo $\Sigma X/n$ por la media de X y extendiendo la ecuación de la fórmula directa:

$$s_x = \sqrt{\frac{\Sigma(X - \bar{X})^2}{n-1}} = \sqrt{\frac{\Sigma[X - (\Sigma X/n)]^2}{n-1}} = \sqrt{\frac{\Sigma X^2 - \frac{(\Sigma X)^2}{n}}{n-1}}$$

Método abreviado para calcular la desviación estándar

$$s_x = \sqrt{\frac{\Sigma X^2 - \frac{(\Sigma X)^2}{n}}{n-1}}$$

donde:

s_x = desviación estándar para la variable de intervalo/razón X

n = tamaño de la muestra

Esta fórmula abreviada no requiere calcular la media ni las puntuaciones de desviación. En la tabla 5-2, simplemente elevamos al cuadrado cada puntuación en bruto en la columna 2, para obtener la columna 5. Entonces las sumas de las columnas 2 y 5 se insertan en la fórmula abreviada:

$$s_x = \sqrt{\frac{\Sigma X^2 - \frac{(\Sigma X)^2}{n}}{n-1}} = \sqrt{\frac{695\,034 - \frac{(2\,856)^2}{12}}{11}} = 37.30 \text{ libras}$$

Una nota precautoria: Tenga cuidado en distinguir entre $(\Sigma X)^2$, la misma suma (de la columna 2) que sirve para calcular la media, y ΣX^2 de la columna 5. Con el propósito de aprender estadística, calcular la desviación estándar mediante ambos métodos resulta una buena manera de detectar los errores de cálculo. Si los dos resultados no son los mismos con un pequeño margen de error de redondeo, verifique todos sus cálculos una vez más.

Limitaciones de la desviación estándar

Ya que la desviación estándar se calcula a partir de la media, al igual que ésta, se infla por los valores extremos. Éstos generan puntuaciones con grandes desviaciones. Cuando se elevan al cuadrado, estas grandes desviaciones, ya sean positivas o negativas, producen una gran puntuación positiva inflada. Así, la desviación estándar puede ser muy confusa cuando se reporta para una distribución sesgada, en la que pocas puntuaciones se extienden en una dirección. Para convencerse del

91

efecto de las puntuaciones extremas tanto en la media como en la desviación estándar, complete la tabla desglosada de la tabla 5-2; pero agregue los dos casos siguientes para obtener una nueva muestra con $n = 14$: el jugador 13 que pesa 115 libras, y el jugador 14 que pesa 125 libras. A continuación compare las respuestas de las muestras original y nueva.

La desviación estándar como parte integral de la estadística inferencial

Las características de la media y de la desviación estándar las hacen muy útiles para alcanzar un sentido de proporción respecto de las variables individuales que se estudian. La desviación estándar y las puntuaciones de desviación, a partir de las cuales se calculan, también son esenciales para examinar las relaciones entre dos variables. El enfoque de la estadística inferencial consiste en desarrollar una comprensión de por qué las puntuaciones individuales de una variable dependiente se desvían de su media.

Suponga, por ejemplo, que estamos estudiando el abuso en el consumo de alcohol. Para una muestra de bebedores adultos, encontramos que la media del consumo de bebidas alcohólicas es de 4.3 galones por año. Gary consumió 7.3 galones el último año, 3 galones arriba de la media. Sam consumió sólo 1 galón, 3.3 galones abajo de la media. ¿Qué sucede con estas desviaciones alta y baja? Quizá podríamos hipotetizar acerca de algunas variables predictoras (independientes) que creamos que estén relacionadas con esta variable dependiente. Por ejemplo, la hipótesis del consumo a la hora de la comida podría explicar, en parte, la puntuación de desviación positiva de Gary: los bebedores de familias que consumen vino con sus alimentos tienen un consumo de alcohol medio más alto. Existe también la hipótesis del bebedor social, la cual podría explicar, en parte, la puntuación de desviación negativa de Sam: los bebedores que sólo consumen alcohol cuando se sirve en convivios sociales tiene un consumo de alcohol medio más bajo.

Para una muestra completa, nuestro interés radica en explicar la variación —la suma de puntuaciones de desviación al cuadrado—. Las puntuaciones de desviación, la variación y la desviación estándar simplemente son medidas de diferencias en las puntuaciones para una variable entre los sujetos de una población. ¿Es más alta la cantidad media de consumo de alcohol anual para las personas de ciertas regiones, entre diferentes edades o grupos religiosos o entre sexos? Las respuestas a tales preguntas dependen de las propiedades matemáticas de la media, la desviación estándar y la curva normal.

¿Por qué se llama desviación "estándar"?

La desviación estándar recibe su nombre del hecho de que proporciona una *unidad de medida común* (un estándar) para comparar variables con *unidades de medida observadas* muy diferentes. Por ejemplo, imagine que Mary Smith y Jason Jones aplican para una beca con base en su desempeño en las pruebas de admisión a la universidad. Mary contestó la prueba académica de la universidad (PAU) y obtuvo 26 puntos PAU. Jason hizo lo propio con la prueba de admisión stanford (PAS) y obtuvo 900 puntos PAS. Estos dos resultados de las pruebas tienen unidades de medida

muy diferentes: los puntos de la prueba PAU varían de cero a 36; y los de la prueba PAS, de 200 a 1 600. Las puntuaciones en bruto para las dos pruebas no pueden compararse directamente. Usando las medias y las desviaciones estándar para ambas pruebas, sin embargo, podemos crear una manera para compararlas. Con los siguientes estadísticos, encontramos que, en comparación con otros aspirantes que contestan estas pruebas, Mary tuvo la puntuación mayor:

X = puntuación de la prueba PAU \bar{X} = 22 puntos PAU s_x = 2 puntos PAU

Y = puntuación de la prueba PAS \bar{Y} = 1 000 puntos PAS s_y = 100 puntos PAS

La puntuación PAU de 26 que obtuvo Mary tiene una desviación estándar de 2 arriba de la media de aquellos que toman la prueba PAU; es decir, su puntuación está 4 puntos PAU —2 veces 2 desviaciones estándar— sobre el promedio de 22. La puntuación de Jason es de 1 *desviación estándar debajo* de la media de aquellos que contestan la prueba PAS; es decir, su puntuación está 100 puntos PAS —1 desviación estándar— debajo del promedio de 1 000. Sin lugar a dudas podemos otorgarle la beca a Mary. Utilizando las desviaciones estándar como unidades de medida en lugar de “puntos de prueba PAU” y de “puntos de prueba PAS”, tenemos una vara de medición común o estándar para ambas variables —de ahí el nombre *desviación estándar*—. ¿Quién le dijo que no podía comparar manzanas con naranjas?

Puntuaciones estandarizadas (puntuaciones Z)

El ejemplo anterior ilustra el hecho de que la puntuación de un sujeto de la investigación en cualquier variable de intervalo/razón puede expresarse de diversas maneras. Primero, lo *expresamos en sus unidades de medida observadas, originales*, como una **puntuación en bruto**. Por ejemplo, la puntuación en bruto X de Mary es 26 puntos PAU. Segundo, lo *expresamos como una desviación de la media* (es decir, la puntuación de desviación $X - \bar{X}$); la puntuación de desviación de Mary es +4 y significa que ella obtuvo 4 puntos PAU arriba de la media de aquellos que tomaron el PAU. Tercero, expresamos su puntuación como un *número de desviaciones estándar de la media* de la puntuación PAU. Llamamos a esto su **puntuación estandarizada** (o puntuación Z), que para la variable X se calcula como sigue:

Cálculo de puntuaciones estandarizadas (puntuaciones Z)

$$Z_x = \frac{X - \bar{X}}{s_x}$$

donde

Z_x = puntuación estandarizada para un valor de X
= número de desviaciones estándar que una puntuación en bruto (puntuación X) se desvía de la media

X = una variable de intervalo/razón

\bar{X} = la media de X

s_x = la desviación estándar de X

¿Por qué se llama desviación "estándar"?

efecto de las puntuaciones extremas tanto en la media como en la desviación estándar, complete la tabla desglosada de la tabla 5-2; pero agregue los dos casos siguientes para obtener una nueva muestra con $n = 14$: el jugador 13 que pesa 115 libras, y el jugador 14 que pesa 125 libras. A continuación compare las respuestas de las muestras original y nueva.

La desviación estándar como parte integral de la estadística inferencial

Las características de la media y de la desviación estándar las hacen muy útiles para alcanzar un sentido de proporción respecto de las variables individuales que se estudian. La desviación estándar y las puntuaciones de desviación, a partir de las cuales se calculan, también son esenciales para examinar las relaciones entre dos variables. El enfoque de la estadística inferencial consiste en desarrollar una comprensión de por qué las puntuaciones individuales de una variable dependiente se desvían de su media.

Suponga, por ejemplo, que estamos estudiando el abuso en el consumo de alcohol. Para una muestra de bebedores adultos, encontramos que la media del consumo de bebidas alcohólicas es de 4.3 galones por año. Gary consumió 7.3 galones el último año, 3 galones arriba de la media. Sam consumió sólo 1 galón, 3.3 galones abajo de la media. ¿Qué sucede con estas desviaciones alta y baja? Quizá podríamos hipotetizar acerca de algunas variables predictoras (independientes) que creamos que estén relacionadas con esta variable dependiente. Por ejemplo, la hipótesis del consumo a la hora de la comida podría explicar, en parte, la puntuación de desviación positiva de Gary: los bebedores de familias que consumen vino con sus alimentos tienen un consumo de alcohol medio más alto. Existe también la hipótesis del bebedor social, la cual podría explicar, en parte, la puntuación de desviación negativa de Sam: los bebedores que sólo consumen alcohol cuando se sirve en convivios sociales tienen un consumo de alcohol medio más bajo.

Para una muestra completa, nuestro interés radica en explicar la variación —la suma de puntuaciones de desviación al cuadrado—. Las puntuaciones de desviación, la variación y la desviación estándar simplemente son medidas de diferencias en las puntuaciones para una variable entre los sujetos de una población. ¿Es más alta la cantidad media de consumo de alcohol anual para las personas de ciertas regiones, entre diferentes edades o grupos religiosos o entre sexos? Las respuestas a tales preguntas dependen de las propiedades matemáticas de la media, la desviación estándar y la curva normal.

¿Por qué se llama desviación "estándar"?

La desviación estándar recibe su nombre del hecho de que proporciona una *unidad de medida común* (un estándar) para comparar variables con *unidades de medida observadas* muy diferentes. Por ejemplo, imagine que Mary Smith y Jason Jones aplican para una beca con base en su desempeño en las pruebas de admisión a la universidad. Mary contestó la prueba académica de la universidad (PAU) y obtuvo 26 puntos PAU. Jason hizo lo propio con la prueba de admisión stanford (PAS) y obtuvo 900 puntos PAS. Estos dos resultados de las pruebas tienen unidades de medida

Si establecemos que X = puntuación PAU con $\bar{X} = 22$ puntos PAU y $s_x = 2$ puntos PAU, la puntuación Z de Mary es

$$Z_x = \frac{X - \bar{X}}{s_x} = \frac{26 - 22}{2} = \frac{4}{2} = 2.00 \text{ DE}$$

donde DE quiere decir "desviaciones estándar". Una puntuación Z es la distancia de una puntuación X hacia la media (es decir, su puntuación de desviación) dividida entre la desviación estándar de las distancias.

Una clave para tener claras estas tres maneras de expresar la puntuación consiste en enfocarse en las unidades de medida. Las puntuaciones en bruto y las puntuaciones de desviación para una variable se presentan en la unidad de medida original observada que, por supuesto, es definida por una variable. Por ejemplo, la unidad de medida observada para edad es años; para peso, libras o kilogramos; para altura, pulgadas o centímetros; y así sucesivamente. Pero no importa de qué unidad de medida de una variable se trate, sus puntuaciones Z se miden en DE. La tabla 5-3 resume estas distinciones.

Aquí aparecen algunos ejemplos de una muestra aleatoria de mujeres estudiantas en la universidad local:

1. Donde X = peso, $\bar{X} = 120$ libras, $s_x = 10$ libras:

Caso	X (peso)	$X - \bar{X}$ (puntuación de desviación)	Z_x (puntuación estandarizada)
Cheryl Jones	110 libras	-10 libras	-1 DE
Jennifer Smith	125 libras	5 libras	.5 DE
Terri Barnett	107 libras	-13 libras	-1.3 DE

2. Donde Y = altura, $\bar{Y} = 65$ pulgadas, $s_y = 3$ pulgadas:

Caso	Y (altura)	$Y - \bar{Y}$ (puntuación de desviación)	Z_y (puntuación estandarizada)
Cheryl Jones	64 pulgadas	-1 pulgada	-.33 DE
Jennifer Smith	65 pulgadas	0 pulgadas	0 DE
Terri Barnett	68 pulgadas	3 pulgadas	1 DE

TABLA 5-3 Las diferentes formas en que se pueden presentar las puntuaciones de una variable

Forma de puntuación para una variable y su símbolo	Unidades de medida de la variable	Ejemplo: $X = \text{altura}$
Puntuación en bruto (puntuación X): X	La unidad de medida de la variable	Pulgadas
Puntuación de desviación = $X - \bar{X}$	La unidad de medida de la variable	Pulgadas
Puntuación estandarizada (Z_x) o "puntuación Z ": $Z_x = \frac{X - \bar{X}}{s_x}$	Desviaciones estándar de la variable (DE)	DE

Tenga presente que tanto las puntuaciones de desviación como las puntuaciones Z son medidas de la distancia desde la puntuación en bruto de una variable hasta su media. La puntuación de desviación se obtiene restando la media de la puntuación en bruto (es decir, $X - \bar{X}$). Al dividir esta puntuación de desviación entre la desviación estándar, cortamos esta puntuación de desviación en las partes y múltiplos de las desviaciones estándar desde la media. Recuerde que después de calcular la media, calcular las puntuaciones de desviación es lo siguiente que hacemos cuando calculamos la desviación estándar. La esencia de la desviación estándar está en ver una puntuación en bruto individual como una desviación desde la media.

Para obtener un buen sentido de proporción sobre las fórmulas para las puntuaciones de desviación y las puntuaciones Z , examinemos las relaciones entre los tamaños de las puntuaciones en bruto, las puntuaciones de desviación y las puntuaciones Z . Primero, cuanto más lejana de la media esté una puntuación X mayores serán su puntuación de desviación y su puntuación Z . Es más, el signo de cualquier puntuación de desviación y puntuación Z indica la dirección de una puntuación: ya sea que la observación caiga arriba de la media (la dirección positiva) o debajo de la media (la dirección negativa). El signo "-" (signo menos) indica que una puntuación en bruto está debajo de la media; el signo "+" (signo más), que está implícito, no escrito, indica que está encima de la media. En los ejemplos anteriores, Cheryl y Terri están debajo del promedio en peso, y Terri está encima del promedio en altura. De hecho, a partir de estas puntuaciones Z podemos decir que Terri es una persona alta, delgada —más de 1 DE debajo en peso, pero 1 DE arriba en altura—. Jennifer tiene altura media; así, su puntuación de desviación y su puntuación Z para Y son cero: ella no se desvía de la altura media.

Puesto que usaremos puntuaciones Z o medidas similares de desviación en cada capítulo en el resto del texto, es prudente practicar cómo calcular puntuaciones de desviación y puntuaciones Z , así como estudiar las direcciones (signos) de esas puntuaciones. Se recomienda una doble verificación. Si una puntuación en bruto queda debajo de la media, su desviación y sus puntuaciones Z son negativas. También tenga presente que las puntuaciones Z son simplemente otra manera de expresar puntuaciones en bruto. Cada puntuación en bruto tiene una puntuación Z correspondiente, y viceversa.

La distribución normal

Además de proporcionar un estándar de comparación entre variables y muestras diferentes, bajo las condiciones apropiadas la media y la desviación estándar ofrecen una riqueza de información. Éste es el caso cuando una variable tiene una distribución de puntuaciones que es normal —formada como la curva de distribución normal—. Como lo definimos en el capítulo 4, una distribución normal es simétrica, con su media, mediana y moda iguales entre sí y localizadas en el centro de la curva. Sin embargo, la simetría o balance en la curva representa la imagen completa. La curva normal también tiene una forma de campana inconfundible, que no es muy plana ni demasiado puntiaguda. Muchas variables se distribuyen

normalmente (como altura, peso e inteligencia). Sin tener en cuenta qué variable se examina, si está normalmente distribuida, posee las propiedades de una curva normal.

Lo que vuelve a la desviación estándar una herramienta estadística tan valiosa es que constituye una parte matemática de la curva normal. Cuando usted sigue la curva desde su centro (es decir, su pico) en cualquier dirección, la curva cambia de forma para acercarse al eje de las X . Desde el pico, el punto en el cual la curva empieza a cambiar hacia afuera es 1 desviación estándar de la media. Este punto se llama punto de inflexión y se destaca en la figura 5-3. Esto indica que la media y la desviación estándar son aspectos matemáticos de un fenómeno natural: la tendencia hacia una distribución normal, en forma de campana para muchos eventos naturales.

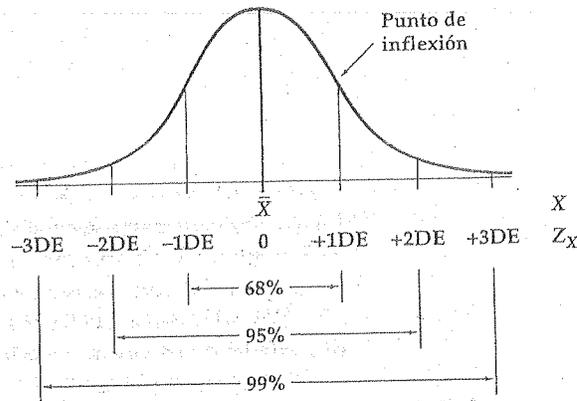
Comprender el fenómeno de normalidad es un aspecto importante de la imaginación estadística. Muchos fenómenos que ocurren naturalmente tienen distribuciones de frecuencias que tienen la forma de campana de la curva normal. La curva normal ilustra el hecho de que cuando nos desviamos más allá de la media, esperamos encontrar cada vez menos casos. Para muchas variables, existe un promedio alrededor del cual cae la mayoría de las puntuaciones, y cuando nos alejamos de este promedio, las frecuencias del caso disminuyen. Por ejemplo, la altura física se distribuye normalmente; la mayoría de las personas está cerca del promedio, con unas cuantas personas muy altas y muy bajas.

Uno de los rasgos más sobresalientes del fenómeno de normalidad que ocurre naturalmente es que ofrece predicciones precisas sobre cuántas puntuaciones de una población caen dentro de cualquier rango de puntuaciones. Como se ilustró en la figura 5-3, *para cualquier variable normalmente distribuida:*

1. Cincuenta por ciento de las puntuaciones caen encima de la media; y 50 por ciento, debajo. Esto se debe al hecho de que la mediana es igual a la media.
2. Virtualmente todas puntuaciones caen dentro de 3 desviaciones estándar a partir de la media en ambas direcciones. Ésta es una distancia de 3 puntuaciones Z debajo hasta 3 puntuaciones Z encima de la media, una amplitud total de 6 desviaciones estándar. La cantidad precisa es 99.7 por ciento. El restante .3 por ciento de casos (es decir, 3 casos de cada 1 000) caen fuera de 3 desviaciones estándar y, teóricamente, la curva se extiende hacia el infinito en ambas direcciones. (Prácticamente hablando, las puntuaciones para algunas variables, como el peso corporal, tienen límites finitos.)
3. Cerca del 95 por ciento de las puntuaciones de una variable normalmente distribuida caen dentro de una distancia de 2 desviaciones estándar, en ambas direcciones de la media. Esto es más menos 2 puntuaciones Z de la media.
4. Aproximadamente 68 por ciento de las puntuaciones de una variable normalmente distribuida caen dentro de una distancia de 1 desviación estándar (más menos 1 puntuación Z), en ambas direcciones de la media.

Tenga presente que la distribución normal tiene características muy predecibles. Si una variable se distribuye en esta peculiar forma de campana, podemos utilizar los estadísticos de la muestra y lo que sabemos respecto de la curva normal, para estimar cuántas puntuaciones en una población caen dentro de cierto rango.

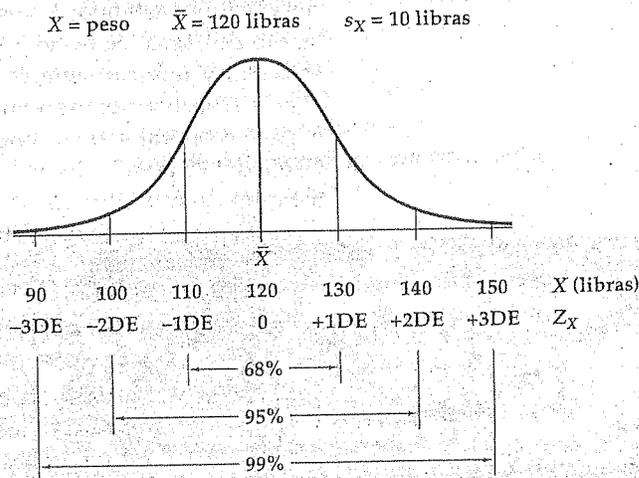
FIGURA 5-3
La relación de la desviación estándar con la curva normal



Para ilustrar la utilidad de la curva normal, sigamos el ejemplo anterior, una muestra de mujeres estudiantes de la universidad local, donde X = peso, el peso medio es de 120 libras, y $s_x = 10$ libras. Primero, necesitamos asegurarnos de que la distribución de puntuaciones es, de hecho, normal; es decir, que tenga forma de campana. Esto podría hacerse produciendo un histograma de puntuaciones de una muestra (no mostrado). Si este gráfico aproximadamente tiene forma de campana, suponemos que esta variable no sólo está normalmente distribuida en la muestra, sino también en la población. Nos referimos a este hecho como "asumiendo la normalidad". (La forma de un histograma de la muestra puede ser ligeramente fuera de lo normal como resultado del error muestral.) Como se grafica en la figura 5-4, asumiendo la normalidad, podemos hacer las siguientes estimaciones de los pesos de la población de mujeres estudiantes en la universidad local:

1. La mitad de estas estudiantes pesa arriba de 120 libras.
2. Cerca del 68 por ciento de las mujeres estudiantes de la universidad local pesan entre 110 y 130 libras.

FIGURA 5-4
Uso de la curva normal para estimar el peso (X) en la distribución de mujeres estudiantes en la universidad local



71

3. Aproximadamente 95 por ciento de las mujeres estudiantes de la universidad local pesan entre 100 y 140 libras.
4. Muy pocas pesan menos de 90 libras o más de 150 libras.

Recuerde, una puntuación Z simplemente es otra forma de expresar una puntuación en bruto (es decir, la puntuación X para una observación individual). Si Susanah pesa 110 libras, ella está 1 DE debajo del peso medio y tiene una puntuación Z de -1.00 DE.

Uso del rango para estimar la desviación estándar

Si una distribución es normal, casi todos los casos caen dentro de la distancia de 3 desviaciones estándar de ambos lados de la media (como se mostró en las figuras 5-3 y 5-4). Esto significa que una variable normalmente distribuida está muy cerca de exactamente 6 desviaciones estándar de amplitud. De hecho, 95 por ciento de las puntuaciones en una distribución normal pueden encontrarse dentro de 2 desviaciones estándar, en ambos lados de la media. Recuerde que el rango proporciona una medida de la expansión de las puntuaciones en una distribución: la distancia del menor al mayor. Si una variable está normalmente distribuida, el tamaño del rango debe ser de aproximadamente 4 a 6 desviaciones estándar de ancho, porque esta dispersión de puntuaciones abarca casi el 100 por ciento de las puntuaciones. Así, para estimar la desviación estándar, dividimos el rango entre 4 o 6, siendo el anterior el método convencional:

Estimación de la desviación estándar usando el rango

$$\text{Estimación } s_x \text{ con base en el rango} = \frac{\text{rango}}{4}$$

donde

s_x = la desviación estándar para la variable X

Rango = (puntuación máxima - puntuación mínima) + el valor de la unidad de redondeo

Antes de calcular la desviación estándar ya sea usando el método directo o el método abreviado, es una buena idea estimarla dividiendo el rango entre 4. Esta estimación y el cálculo subsecuente de la desviación estándar deben tener valores cercanos. Si no lo son, es una señal de que la distribución está sesgada o que los valores extremos están incrementando el valor del rango. En estas situaciones, la observación de un histograma arroja luz acerca de la diferencia entre los valores estimados y los valores calculados.

Con sus vínculos matemáticos con la media y la desviación estándar, la curva normal es un recurso muy útil, particularmente como una distribución de probabilidad, el tema principal del capítulo 6.

Una ilustración completa del cálculo de los estadísticos de dispersión

Ahora que hemos discutido los conceptos que hay detrás del rango y de la desviación estándar, completemos un problema para ilustrar la organización básica de la estadística descriptiva. Los tres estadísticos más comúnmente reportados para datos de intervalo/razón son la media, la desviación estándar y el rango.

Uso de un formato desglosado para calcular la desviación estándar

La tabla 5-4 presenta una tabla desglosada de los impuestos a la gasolina cobrados por estados del oeste seleccionados; así, X = impuesto de gasolina por galón, y hemos ordenado las puntuaciones.

Primero, calculemos el rango. Con estas puntuaciones ordenadas, vemos que la puntuación mínima es 17 centavos; y la máxima 28 centavos. Nuestra unidad de redondeo es un número entero.

$$\text{Rango} = (\text{puntuación máxima} - \text{puntuación mínima}) + \text{valor de la unidad de redondeo} = (28 - 17) + 1 = 12¢$$

Segundo, estimamos la desviación estándar utilizando el rango:

$$\text{Estimación de } s_x \text{ basada en el rango} = \frac{\text{rango}}{4} = \frac{12}{4} = 3¢$$

Tercero, calculamos la media y la usamos para calcular las puntuaciones de desviación y para completar las sumas en la tabla desglosada (tabla 5-4):

TABLA 5-4 Impuestos estatales en la gasolina en estados del oeste seleccionados durante mayo de 1996

Especificaciones		Cálculos		
Estado	Impuesto (¢) por galón X	Desviaciones $X - \bar{X}$	$(X - \bar{X})^2$	X^2
Nuevo México	17	-4.7	22.09	289
California	18	-3.7	13.69	324
Arizona	18	-3.7	13.69	324
Utah	19	-2.7	7.29	361
Colorado	22	.3	.09	484
Washington	23	1.3	1.69	529
Nevada	23	1.3	1.69	529
Oregon	24	2.3	5.29	576
Idaho	25	3.3	10.89	625
Montana	28	6.3	39.69	784
$\Sigma X = 217 ¢$		$\Sigma(X - \bar{X})^2 = 116.10 ¢ \text{ cuadrados}$		
$n = 10$		$\Sigma(X - \bar{X}) = 0$		
		$\Sigma X^2 = 4 825 ¢ \text{ cuadrados}$		

FUENTE: Tasas de impuesto de <http://www.api.org/news/596sttax.htm>. Copyright © 1996 por el American Petroleum Institute. Reimpreso con autorización del Instituto.

$$\bar{X} = \frac{\Sigma X}{n} = \frac{217}{10} = 21.7¢$$

Cuarto, calculamos la desviación estándar con el método directo:

$$s_x = \sqrt{\frac{\Sigma(X - \bar{X})^2}{n - 1}} = \sqrt{\frac{116.10}{9}} = 3.59¢$$

Quinto, calculamos la desviación estándar usando el método abreviado:

$$s_x = \sqrt{\frac{\Sigma X^2 - \frac{(\Sigma X)^2}{n}}{n - 1}} = \sqrt{\frac{4825 - \frac{(217)^2}{10}}{9}} = 3.59¢$$

Sexto, verificamos si ambos cálculos de desviación estándar produjeron la misma respuesta. Éste es el caso. Por último, comparamos la desviación estándar calculada con la estimación y observamos que están relativamente cerca. Si la estimación hubiera sido de la mitad del tamaño o el doble del valor calculado, eso habría justificado analizar la existencia de valores extremos u otras peculiaridades en la distribución de las puntuaciones. Ahora que tenemos la media y la desviación estándar, vemos que el promedio de impuestos a la gasolina promedia cerca de 22 centavos por galón y que aproximadamente dos de cada tres estados (aproximadamente 68 por ciento) están dentro de 3.59 centavos de este promedio.

Uso de un formato de distribución de frecuencias para calcular la desviación estándar

En el capítulo 4 vimos que un formato de distribución de frecuencias es una forma más concisa de organizar datos. Usar este formato simplemente requiere contar cada puntuación el número de veces f que ocurre. La tabla 5-5 presenta los datos sobre los impuestos a la gasolina como una distribución de frecuencias. Calculemos la desviación estándar con fórmulas modificadas para considerar la frecuencia (f) de cada puntuación.

Primero, calculemos el rango. En la tabla 5-5 las puntuaciones se ordenan con una puntuación mínima de 17 centavos y una máxima de 28 centavos. Nuestra unidad de redondeo es un número entero. Así,

$$\begin{aligned} \text{Rango} &= (\text{puntuación máxima} - \text{puntuación mínima}) + \text{valor de la unidad} \\ &\quad \text{de redondeo} \\ &= (28 - 17) + 1 = 12¢ \end{aligned}$$

Segundo, estimamos la desviación estándar usando el rango, como lo hicimos anteriormente:

$$\text{Estimación } s_x \text{ basada en el rango} = \frac{\text{rango}}{4} = \frac{12}{4} = 3¢$$

Tercero, calculamos la media y la usamos para calcular las puntuaciones de desviación y para completar las sumas en la tabla 5-5:

$$\bar{X} = \frac{\Sigma f(X)}{n} = \frac{217}{10} = 21.7¢$$

TABLA 5-5 Cálculo de la desviación estándar usando un formato de distribución de frecuencias
(Impuestos estatales en la gasolina para los estados del oeste seleccionados durante mayo de 1996)

Especificaciones		Cálculos							
X	f	f(X)	(X - \bar{X})	f(X - \bar{X})	(X - \bar{X}) ²	f(X - \bar{X}) ²	X ²	f(X ²)	
17	1	17	-4.7	-4.7	22.09	22.09	289	289	
18	2	36	-3.7	-7.4	13.69	27.38	324	648	
19	1	19	-2.7	-2.7	7.29	7.29	361	361	
22	1	22	.3	.3	.09	.09	484	484	
23	2	46	1.3	2.6	1.69	3.38	529	1058	
24	1	24	2.3	2.3	5.29	5.29	576	576	
25	1	25	3.3	3.3	10.89	10.89	625	625	
28	1	28	6.3	6.3	39.69	39.69	784	784	
		$\Sigma f(X) = 217¢$				$\Sigma f(X - \bar{X})^2 = 116.10¢$ cuadrados			
		n = 10	$\Sigma f(X - \bar{X}) = 0$			$\Sigma f(X^2) = 4 825¢$ cuadrados			

FUENTE: Tasas de impuesto de <http://www.api.org/news/596sttax.htm>. Copyright © 1996 por el American Petroleum Institute. Reimpreso con autorización del Instituto.

Cuarto, calculamos la desviación estándar con el método directo:

$$s_x = \sqrt{\frac{\Sigma f(X - \bar{X})^2}{n - 1}} = \sqrt{\frac{116.10}{9}} = 3.59¢$$

Quinto, calculamos la desviación estándar usando el método abreviado:

$$s_x = \sqrt{\frac{\Sigma fX^2 - \frac{[\Sigma f(X)]^2}{n}}{n - 1}} = \sqrt{\frac{4 825 - \frac{(217)^2}{10}}{9}} = 3.59¢$$

Sexto, verificamos si los dos cálculos son iguales, lo cual es afirmativo. Por último, comparamos la desviación estándar calculada con la estimación y observamos que están relativamente cerca.

Presentación tabular de resultados

En artículos de investigación, una tabla básica de estadística descriptiva es la que lista todas las variables y sus medias y desviaciones estándar. La tabla 5-6 presenta una tabla de estadística descriptiva de un estudio del bienestar psicológico de personas sin hogar en dos puntos en el tiempo.

TABLA 5-7 La distribución sesgada de estancias en el hospital durante el último año entre personas mayores de 65 años de edad (datos ficticios)

Especificaciones		Cálculos		
(1) Caso	(2) X	(3) $X - \bar{X}$	(4) $(X - \bar{X})^2$	(5) X^2
1	0	-2.41	5.81	0
2	0	-2.41	5.81	0
3	0	-2.41	5.81	0
4	0	-2.41	5.81	0
5	0	-2.41	5.81	0
6	0	-2.41	5.81	0
7	0	-2.41	5.81	0
8	0	-2.41	5.81	0
9	1	-1.41	1.99	1
10	1	-1.41	1.99	1
11	1	-1.41	1.99	1
12	2	-.41	.17	4
13	2	-.41	.17	4
14	5	2.59	6.71	25
15	9	6.59	43.43	81
16	10	7.59	57.61	100
17	10	7.59	57.61	100
$\Sigma X = 41$ veces		$\Sigma(X - \bar{X})^2 = 218.15$ veces		
$n = 17$	$\Sigma(X - \bar{X}) = .03^*$	$\Sigma X^2 = 317$ veces cuadradas		

* No sumó cero debido al error de redondeo.

Incluso sin un histograma, los valores relativos de la media y de la desviación estándar para esta distribución proporcionan una señal de que la distribución está sesgada. Estos estadísticos se calculan como sigue:

X = estancias en el hospital = durante el último año, el número de veces que una persona es admitida en un hospital y pasa por lo menos una noche

$$\bar{X} = 2.41 \text{ veces} \quad s_x = 3.69 \text{ veces} \quad n = 17 \text{ casos}$$

Observe que la desviación estándar es más grande que la media. Esto sugiere que una o más puntuaciones extremas inflaron la media y la desviación estándar. Por otra parte, desde el momento en que se elevan al cuadrado los números en la desviación estándar, unas cuantas puntuaciones extremas pueden hacer "explotar" rápidamente su valor. Note, por ejemplo, la enorme contribución a la suma de cuadrados que los tres casos más grandes hicieron con sus estancias de 9, 10 y 10 veces.

¿Por qué una desviación estándar más grande que la media indica un sesgo? Recuerde que si una distribución no está sesgada (es decir, tiene una forma de campana normal); su rango tendrá una amplitud de aproximadamente 4 a 6 desviaciones estándar. Cuando la curva es trazada, la amplitud de 2 o 3 desviaciones estándar se ajustará en cada lado de la media. Si el límite inferior de las puntuaciones X de una variable es cero, por lo menos la distancia de 2 desviaciones estándar debería ajustarse entre una puntuación X de cero y la media. Cuando la desviación estándar es más grande que la media, como en el caso de las estancias en el hospital, ni una sola amplitud de la desviación estándar puede lograr este ajuste. Otra manera de explicarlo es que la desviación estándar debería ser alrededor de la mitad del tamaño de la media o menos.

Dos reglas generales se aplican a los tamaños relativos de la media y de la desviación estándar:

1. Si la desviación estándar es más grande que la media, esto probablemente indica un sesgo, es decir, la presencia de valores extremos u otra peculiaridad en la forma de la distribución, como una distribución bimodal.
2. Si la desviación estándar no es de la mitad de tamaño de la media o menos, se debe tener cuidado al examinar la distribución para analizar la posible existencia de sesgos o valores extremos.

Como veremos en capítulos posteriores, cuando una variable sesgada es correlacionada con otras variables, los resultados pueden ser erróneos (capítulo 14). En tales casos, deben realizarse ajustes a los estadísticos para evitar tales errores.

Fórmulas en el capítulo 5

Organice una tabla desglosada con los casos ordenados de acuerdo al rango:

Especificaciones		Cálculos		
(1) Caso	(2) X	(3) $X - \bar{X}$	(4) $(X - \bar{X})^2$	(5) X^2
•	•
•	•
•	•
$\Sigma X = \dots$		$\Sigma(X - \bar{X})^2 = \dots$		
$n = \dots$	$\Sigma(X - \bar{X}) = 0$		$\Sigma X^2 = \dots$	

TABLA 5-6 Estadísticos descriptivos para los síntomas psicológicos, satisfacción de vida y manejo personal

Subescalas	Seguimiento 1		Seguimiento 2	
	M	DE	M	DE
Síntomas psicológicos				
Enojo	4.17	.80	4.14	.85
Ansiedad	3.97	.79	3.97	.80
Depresión	3.60	.76	3.68	.77
Manía	3.59	.87	3.68	.90
Psicosis	4.51	.72	4.52	.72
Satisfacción de vida				
Vestido	4.33	1.59	4.49	1.60
Alimentación	4.79	1.53	4.98	1.42
Salud	4.81	1.38	4.77	1.41
Vivienda	4.37	1.49	4.51	1.54
Diversión	3.74	1.53	3.84	1.56
Dinero	2.98	1.57	3.19	1.67
Social	4.42	1.44	4.51	1.79
Manejo personal				
Manejo-1	3.21	.85	3.24	.84
Manejo-2	3.36	.87	3.28	.85

NOTA: $n = 298$. Las puntuaciones altas reflejan mayor bienestar subjetivo.

FUENTE: Modificada de Marshall y *et al.*, 1996: 49.

⊗ INSENSATEZ Y FALACIAS ESTADÍSTICAS ⊗

¿Qué indica cuando la desviación estándar es más grande que la media?

Como vimos en el capítulo 4, la media es susceptible de distorsión por la presencia de puntuaciones extremas, valores extremos y distribuciones sesgadas. Como se basa en desviaciones de la media, la desviación estándar es susceptible del mismo problema. La distorsión está determinada por el hecho de que las puntuaciones de desviación están elevadas al cuadrado.

Un tipo común de distribución sesgada es un sesgo positivo (o derecho), en el cual la mayoría de las personas tienen bajas puntuaciones, pero algunas obtienen altas puntuaciones. Por ejemplo, "la estancia en el hospital", o el número de veces que una muestra aleatoria de personas mayores de 65 años han permanecido en un hospital durante el último año, es un sesgo derecho. La mayoría de las personas registrará cero en estancia; algunas, uno; otras reportarán dos, y pocas personas muy enfermas, anotarán estancias frecuentes. Este tipo de distribución se presenta en la tabla 5-7.

O si lo desea, organice los datos en una distribución de frecuencias con los casos ordenados de acuerdo al rango:

Especificaciones		Cálculos						
X	f	f(X)	(X - \bar{X})	f(X - \bar{X})	(X - \bar{X}) ²	f(X - \bar{X}) ²	X ²	f(X ²)
•	•
•	•
•	•
n = ...		$\Sigma f(X) = \dots$	$\Sigma f(X - \bar{X}) = 0$		$\Sigma f(X - \bar{X})^2 = \dots$	$\Sigma f(X^2) = \dots$		

Para calcular el rango:

1. Ordene las puntuaciones en la distribución de menor a mayor.
2. Identifique las puntuaciones mínima y máxima.
3. Identifique el valor de la unidad de redondeo (véase apéndice A).
4. Calcule el rango:

$$\text{Rango} = (\text{puntuación máxima} - \text{puntuación mínima}) + \text{valor de la unidad de redondeo}$$

Estimación de la desviación estándar usando el rango:

$$\text{Estimación de } s_x \text{ basada en el rango} = \frac{\text{Rango}}{4}$$

Método directo para calcular la desviación estándar:

1. Empiece calculando la media de X y completando una tabla desglosada similar a la de la tabla 5-2.
2. Calcule la desviación estándar:

Mediante una tabla desglosada Mediante una distribución de frecuencias

$$s_x = \sqrt{\frac{\Sigma(X - \bar{X})^2}{n - 1}}$$

$$s_x = \sqrt{\frac{\Sigma f(X - \bar{X})^2}{n - 1}}$$

Método abreviado para calcular la desviación estándar:

Mediante una tabla desglosada Mediante una distribución de frecuencias

$$s_x = \sqrt{\frac{\Sigma X^2 - \frac{(\Sigma X)^2}{n}}{n - 1}}$$

$$s_x = \sqrt{\frac{\Sigma fX^2 - \frac{[\Sigma f(X)]^2}{n}}{n - 1}}$$

Cálculo de puntuaciones estandarizadas (puntuaciones Z):

$$Z_x = \frac{X - \bar{X}}{s_x}$$

Preguntas para el capítulo 5

1. Los estadísticos de dispersión se calculan sólo en algunos tipos de variables. ¿Qué niveles de medición de las variables admiten los anteriores cálculos?
2. Tanto el rango como la desviación estándar son medidas de la dispersión de las puntuaciones en una distribución. Explique las diferencias en perspectiva entre estos dos estadísticos.
3. ¿Qué efecto tiene una puntuación extrema o valor extremo sobre el cálculo del rango?
4. La desviación estándar se "deriva" de la media. ¿Qué significa esto?
5. Al calcular el rango, el valor de la unidad de redondeo de la variable se suma a la diferencia entre las puntuaciones máxima y mínima. ¿Por qué se suma el valor de la unidad de redondeo?
6. Al calcular la desviación estándar, ¿por qué es necesario elevar al cuadrado las puntuaciones de desviación?
7. Al calcular la desviación estándar para datos de una muestra, ¿por qué debemos dividir entre $n - 1$?
8. Al calcular la desviación estándar, ¿por qué se requiere tomar la raíz cuadrada?
9. ¿Cuál es la relación matemática entre la varianza y la desviación estándar?
10. Mencione otro nombre para la variación.
11. ¿Cuál es el significado de la palabra *estándar* en el término *desviación estándar*?
12. Una expresión de qué tan lejos está una puntuación en bruto de la media de una distribución, en las unidades de medida originales de la variable X, se llama una puntuación _____.
13. Una expresión de qué tan lejos está una puntuación en bruto de la media de una distribución, en unidades de medida de desviaciones estándar (DE), se llama una puntuación _____.
14. ¿Cuáles son las propiedades de una distribución normal?
15. En una distribución normal, ¿qué porcentaje de puntuaciones caen aproximadamente dentro de 1 desviación estándar de la media en ambas direcciones?, ¿y dentro de 2 desviaciones estándar de la media en ambas direcciones?, ¿y dentro de 3 desviaciones estándar de la media en ambas direcciones?
16. En una distribución normal, ¿qué porcentaje de puntuaciones caen exactamente sobre la media? ¿Qué estadístico de tendencia central, además de la media, cuenta para este fenómeno?
17. En una distribución normal, la curva alcanza su máximo en el valor de la media. ¿Qué estadístico de tendencia central, además de la media, justifica este fenómeno?

paquete, y sólo tres son ases. Debe prestarse mucha atención en las cuestiones de reemplazamiento en eventos compuestos. Se ajustan numeradores y denominadores consecuentemente. Por ejemplo, calculemos lo siguiente:

$$p [\text{as luego rey luego as}] \text{ sin reemplazamiento} = p [\text{as}] \cdot p [\text{rey}] \cdot p [\text{as}]$$

$$= \frac{4}{52} \cdot \frac{4}{51} \cdot \frac{3}{50} = \frac{48}{132\,600} = .0004$$

Finalmente, no todos los eventos compuestos implican cuestiones de reemplazamiento. Por ejemplo, el reemplazamiento no es una cuestión que tenga que ver en el lanzamiento de una moneda. Las probabilidades calculadas son las mismas para "cara luego cara" al lanzar dos monedas a la vez o al lanzar una moneda dos veces.

Las cinco reglas de probabilidad son fundamentales; es decir, deben considerarse para calcular la probabilidad de cualquier evento, no importa lo simple o complicado que el evento sea. Los ejemplos simples presentados en este capítulo ilustran tales principios básicos. Formulaciones mucho más complejas de probabilidades se presentan en textos avanzados como el de Lee y Maykovich (1995). Por fortuna, para los estudiantes y los investigadores de hoy, no es necesario tener habilidades matemáticas extensas para calcular probabilidades. Los programas de computación sólo requieren que aprendamos qué botones oprimir o qué tabla leer para obtener respuestas a las preguntas de probabilidad. Sin embargo, una comprensión cabal de la teoría básica de probabilidad es necesaria para evitar interpretaciones erróneas de dicha información de la computadora. Es más, una comprensión de la teoría de la probabilidad es esencial para adquirir la imaginación estadística.

Uso de la curva normal como una distribución de probabilidad

Pensamiento proporcional respecto de un grupo de casos y casos únicos

Como observamos en el capítulo 5, la desviación estándar sirve para examinar la forma en que las puntuaciones se dispersan en una distribución, y para comparar la dispersión de dos o más muestras. Sin embargo, podemos lograr mucho más con la desviación estándar. Con una sola variable de intervalo/razón *que tenemos razón para creer que está normalmente distribuida en su población*, podemos calcular puntuaciones estandarizadas (puntuaciones Z) y usarlas para determinar la proporción (p) de puntuaciones de una población que caen entre cualesquiera dos puntuaciones en la distribución. Puesto que la curva normal tiene una forma inconfundible, podemos identificar y medir áreas porque éstas representan una proporción de casos.

Recuerde del capítulo 5 que una puntuación Z nos dice cuántas desviaciones estándar está alejada una puntuación bruta (o puntuación X) de la media:

$$Z_x = \frac{X - \bar{X}}{s_x} = \text{número de desviaciones estándar (DE) desde la media}$$