

# Descripción de datos con medidas numéricas

## Estadística



Gamaliel Moreno Chávez

Centro de Crecimiento Humanista

Enero  
2021



# Generalidades

---

Las gráficas pueden ayudar a describir la forma básica de una distribución de datos pero hay limitaciones

- No siempre es posible presentar gráficas.
- Son poco precisas.
- No son comparables

Una forma de superar estos problemas es usar medidas numéricas, que se pueden calcular para una muestra o una población de mediciones.

# Generalidades

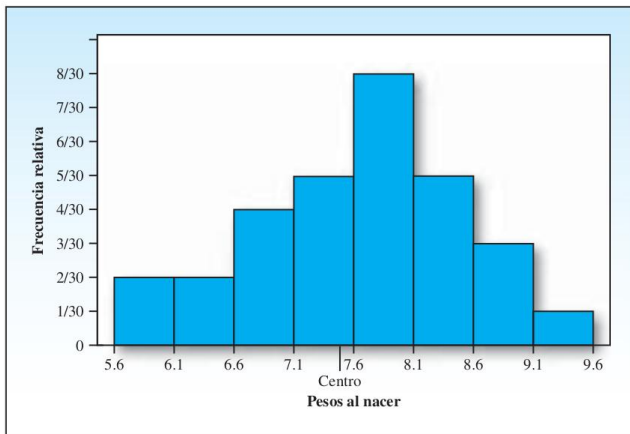
---

## Definición

Las mediciones descriptivas numéricas asociadas con una población de mediciones se llaman parámetros; las calculadas a partir de mediciones muestrales reciben el nombre de estadísticas.

## Medidas de centro

Una de las primeras mediciones numéricas importantes es una medida de centro, es decir, una medida a lo largo del eje horizontal que localiza el centro de la distribución.



# Media

---

El **promedio aritmético** de un conjunto de mediciones es una medida de centro muy común y útil. Se conoce como **media aritmética** o simplemente **media**. Para distinguirla usamos el símbolo  $\bar{x}$  (x barra).

## Definición

La media aritmética o promedio de un conjunto de  $n$  mediciones es igual a la suma de las mediciones dividida entre  $n$

# Media

---

## NOTACIÓN

Media muestral:  $\bar{x} = \frac{\sum x_i}{n}$

Media poblacional:  $\mu$

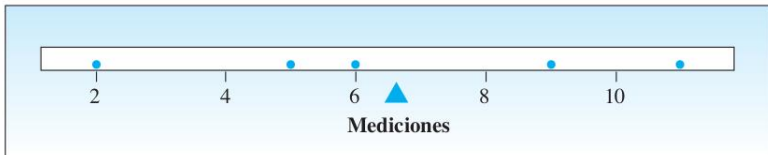
$\sum_{i=1}^n x_i$  que significa  $x_1 + x_2 + x_3 + \cdots + x_n$

$\sum x_i$  que significa “la suma de todas las mediciones de x”

# Media

---

## Ejemplo



$$\bar{x} = \frac{\sum x_i}{n} = \frac{2 + 9 + 11 + 5 + 6}{5} = 6.6$$

# Media

---

Ejercicio. Calcule la media para los siguientes datos

## **Pesos de 30 bebés de gestación**

7.2	7.8	6.8	6.2	8.2
8.0	8.2	5.6	8.6	7.1
8.2	7.7	7.5	7.2	7.7
5.8	6.8	6.8	8.5	7.5
6.1	7.9	9.4	9.0	7.8
8.5	9.0	7.7	6.7	7.7



# Mediana

---

Una segunda medida de tendencia central es la mediana, que es el valor de la posición media en el conjunto de mediciones ordenada de menor a mayor.

## Definición

La mediana  $m$  de un conjunto de  $n$  mediciones es el valor de  $x$  que cae en la posición media cuando las mediciones son ordenadas de menor a mayor.

# Mediana

---

## Ejemplo

Encuentre la mediana para el conjunto de mediciones 2, 9, 11, 5, 6.

**Solución** Ordene las  $n = 5$  mediciones de menor a mayor:

2   5   6   9   11  
          ↑

La observación de enmedio, marcada con una flecha, es el centro del conjunto o sea  $m = 6$ .

# Mediana

---

## Ejemplo

Encuentre la mediana para el conjunto de mediciones 2, 9, 11, 5, 6, 27.

**Solución** Ordene las mediciones de menor a mayor:

2   5   6   9   11   27  
          ↑

Ahora hay dos observaciones “de enmedio”, vistas en la caja. Para hallar la mediana, escoja un valor a la mitad entre las dos observaciones de enmedio:

$$m = \frac{6 + 9}{2} = 7.5$$

# Mediana

---

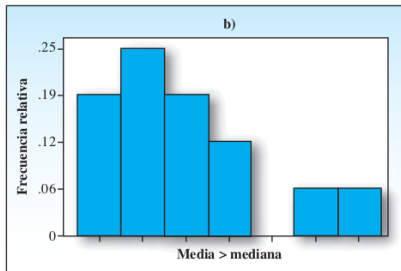
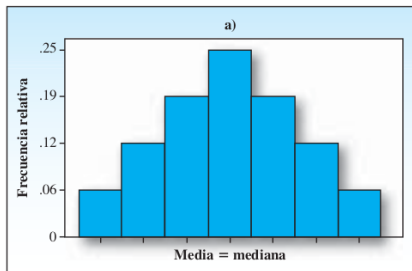
Aunque tanto la media como la mediana son buenas medidas del centro de una distribución, la mediana es menos sensible a valores extremos o resultados atípicos. Por ejemplo, el valor  $x = 27$  en el ejemplo anterior es mucho mayor que las otras mediciones. La mediana,  $m = 7.5$ , no es afectada por el resultado atípico, en tanto que el promedio muestral,

$$\bar{x} = \frac{\sum x_i}{n} = \frac{60}{6} = 10$$

sí es afectado.

# Mediana

Cuando un conjunto de datos tiene valores extremadamente pequeños u observaciones muy grandes, la media muestral se traza hacia la dirección de las mediciones extremas.



# Moda

---

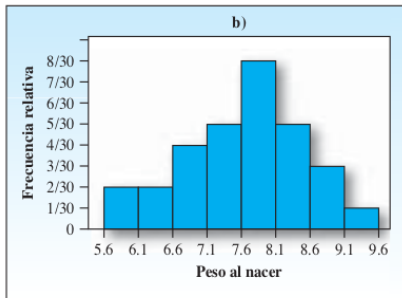
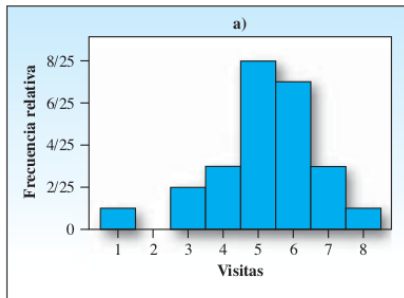
Otra forma de localizar el centro de una distribución es buscar el valor de  $x$  que se presenta con la frecuencia más alta. Esta medida del centro se denomina moda.

## Definición

La moda es la categoría que se presenta con más frecuencia o el valor de  $x$  que se presenta con más frecuencia. Cuando las mediciones en una variable continua se han agrupado como histograma de frecuencia o de frecuencia relativa, la clase con el pico más alto o frecuencia se llama clase modal, y el punto medio de esa clase se toma como la moda.



# Moda



# Moda

---

Es posible que una distribución de mediciones tenga más de una moda. Estas modas aparecerían como “picos locales” en la distribución de frecuencia relativa. Por ejemplo, si fuéramos a tabular la longitud de los peces sacados de un lago durante una temporada, podríamos obtener una distribución bimodal, posiblemente reflejando una mezcla de peces jóvenes y viejos en la población. A veces las distribuciones bimodales de tamaños o pesos reflejan una mezcla de mediciones tomadas en machos y hembras. En cualquier caso, un conjunto o distribución de mediciones puede tener más de una moda.



## Ejercicio

Tiempo transcurrido entre la toma de pedido y el servicio a la mesa de los restaurantes fueron registrados

Tiempo	Restaurante A									
	32	12	24	24	23	21	27	23	23	23
	12	17	21	11	19	14	20	15	23	24
	10	25	18	17	21	21	13	21	21	19
	Restaurante B									
	14	22	14	16	13	11	7	13	13	16
	22	22	11	21	25	24	12	25	13	16
	20	17	28	27	11	21	23	11	11	25

¿Cuál restaurante considera usted que presta un mejor servicio y por qué?

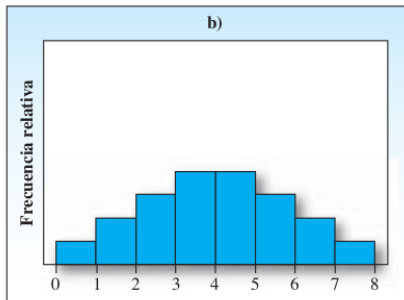
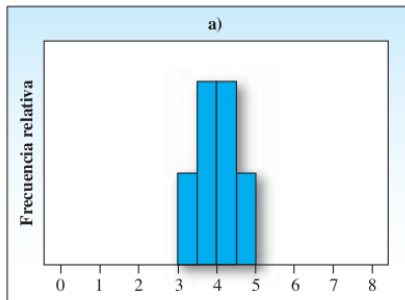
# Ejercicio

## Resultados

Descriptive statistics		Descriptive statistics	
<i>Restaurante A</i>		<i>Restaurante B</i>	
count	30	count	30
mean	19.77	mean	17.47
1st quartile	17.00	1st quartile	13.00
median	21.00	median	16.00
3rd quartile	23.00	3rd quartile	22.00
interquartile r:	6.00	interquartile r:	9.00
mode	21.00	mode	11.00
low extremes	0	low extremes	0
low outliers	0	low outliers	0
high outliers	0	high outliers	0
high extreme:	0	high extreme:	0

## Medidas de variabilidad

Los conjuntos de datos pueden tener el mismo centro pero con aspecto diferente por la forma en que los números se dispersan desde el centro. Considere las dos distribuciones que se muestran en la figura 2.6. Ambas distribuciones están centradas en  $x = 4$ , pero hay una gran diferencia en la forma en que las mediciones se dispersan o varían.



# Rango

---

## Definición

El rango,  $R$ , de un conjunto de  $n$  mediciones se define como la diferencia entre la medición más grande y la más pequeña.

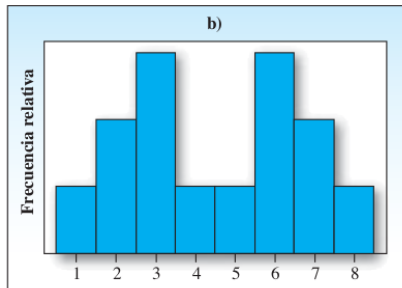
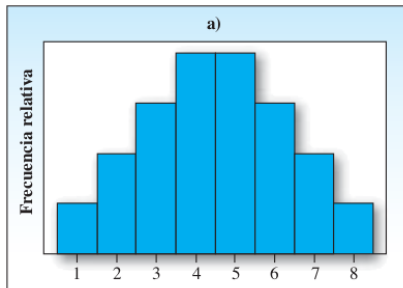
Para los datos de peso al nacer de la tabla, las mediciones varían de 5.6 a 9.4. Por tanto, el rango es  $9.4 - 5.6 = 3.8$ .

### Pesos de 30 bebés de gestación

7.2	7.8	6.8	6.2	8.2
8.0	8.2	5.6	8.6	7.1
8.2	7.7	7.5	7.2	7.7
5.8	6.8	6.8	8.5	7.5
6.1	7.9	9.4	9.0	7.8
8.5	9.0	7.7	6.7	7.7

# Rango

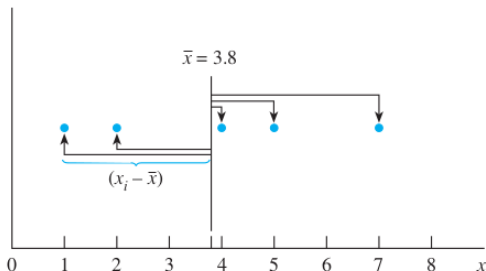
El rango es fácil de calcular, fácil de interpretar y es una medida adecuada de variación para conjuntos pequeños de datos. Pero, para conjuntos grandes, el rango no es una medida adecuada de variabilidad. Por ejemplo, las dos distribuciones de frecuencia relativa de la figura tienen el mismo rango pero muy diferentes formas y variabilidad.



# Varianza

¿Hay una medida de variabilidad que sea más sensible que el rango?

$$\bar{x} = \frac{\sum x_i}{n} = \frac{19}{5} = 3.8$$



## Cálculo de $\sum (x_i - \bar{x})^2$

$x$	$(x_i - \bar{x})$	$(x_i - \bar{x})^2$
5	1.2	1.44
7	3.2	10.24
1	-2.8	7.84
2	-1.8	3.24
4	.2	.04
19	0.0	22.80

# Varianza

---

De la suma de desviaciones cuadradas, se calcula una sola medida llamada varianza. Para la varianza de una muestra usamos el símbolo  $s^2$  y la varianza de una población  $\sigma^2$ . La varianza será relativamente grande para datos muy variables y relativamente pequeña para datos menos variables.

## Definición

La **varianza de una población** de  $N$  mediciones es el promedio de los cuadrados de las desviaciones de las mediciones alrededor de su media  $\mu$ . La varianza poblacional se denota con  $\sigma^2$  y está dada por la fórmula

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$$

# Varianza

---

## Definición

La **varianza de una muestra** de  $n$  mediciones es la suma de las desviaciones cuadradas de las mediciones alrededor la media  $\bar{x}$  dividida entre  $(n - 1)$ . La varianza muestral se denota con  $s^2$  y está dada por la fórmula

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



# Varianza

Para el conjunto de  $n = 5$  mediciones muestrales presentadas en la tabla

suma

$$\sum (x_i - \bar{x})^2 = 22.80$$

varianza muestral

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1} = \frac{22.80}{4} = 5.70$$

## Cálculo de $\sum (x_i - \bar{x})^2$

$x$	$(x_i - \bar{x})$	$(x_i - \bar{x})^2$
5	1.2	1.44
7	3.2	10.24
1	-2.8	7.84
2	-1.8	3.24
4	.2	.04
19	0.0	22.80

# Desviación estándar

---

La varianza se mide en términos del cuadrado de las unidades originales de medición. Tomando la raíz cuadrada de la varianza, obtenemos la desviación estándar, que regresa la medida de variabilidad a las unidades originales de medición.

## Definición

La **desviación estándar** de un conjunto de mediciones es igual a la raíz cuadrada positiva de la varianza.

# Notación

---

## NOTACIÓN

$n$ : número de mediciones en la muestra

$s^2$ : varianza muestral

$s = \sqrt{s^2}$ : desviación muestral estándar

$N$ : número de mediciones en la población

$\sigma^2$ : varianza poblacional

$\sigma = \sqrt{\sigma^2}$ : desviación poblacional estándar

# Resumen

---

- El valor de  $s$  es siempre mayor o igual a cero.
- Cuanto mayor sea el valor de  $s^2$  o de  $s$ , mayor es la variabilidad del conjunto de datos.
- Si  $s^2$  o  $s$  es igual a cero, todas las mediciones deben tener el mismo valor.
- Para medir la variabilidad en las mismas unidades que las observaciones originales, calculamos la desviación estándar  $s = \sqrt{s^2}$ .

## Coeficiente de variación

---

El coeficiente de variación es la relación entre la desviación típica de una muestra y su media.

$$CV = \left( \frac{\sigma}{\bar{x}} \right) 100$$

El coeficiente de variación permite comparar las dispersiones de dos distribuciones distintas, siempre que sus medias sean positivas. Mayor dispersión el valor del coeficiente de variación será mayor.

## Coeficiente de variación

---

Una distribución tiene  $\bar{x} = 140$  y  $\sigma = 28.28$  y otra  $\bar{x} = 150$  y  $\sigma = 24$ .  
¿Cuál de las dos presenta mayor dispersión?

$$CV_1 = \frac{28,28}{140} \cdot 100 = 20,2 \%$$

$$CV_2 = \frac{24}{150} \cdot 100 = 16 \%$$

La primera distribución presenta mayor dispersión.

# Sobre la significancia práctica de la desviación estándar

---

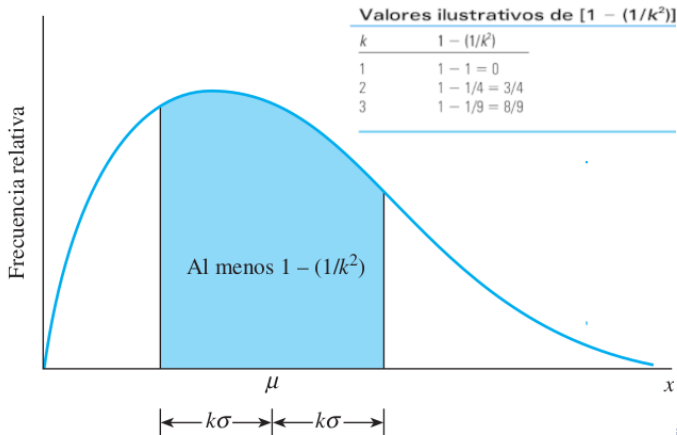
## Teorema de Chebyshev

Dado un número  $k$  mayor o igual a 1 y un conjunto de  $n$  mediciones, al menos  $[1 - (1/k^2)]$  de las mediciones estarán dentro de  $k$  desviaciones estándar de su media.

Se construye un intervalo al medir una distancia  $k\sigma$  a cualquier lado de la media  $m$ . El número  $k$  puede ser cualquier número mientras sea mayor o igual a 1. Entonces el teorema de Chebyshev expresa que al menos  $[1 - (1/k^2)]$  del número total  $n$  de mediciones está en el intervalo construido.

# Teorema de Chebyshev

## Ilustración del teorema de Chebyshev





# Teorema de Chebyshev

---

## Ilustración del teorema de Chebyshev

### Valores ilustrativos de $[1 - (1/k^2)]$

$k$	$1 - (1/k^2)$
1	$1 - 1 = 0$
2	$1 - 1/4 = 3/4$
3	$1 - 1/9 = 8/9$

- Al menos ninguna de las mediciones está en el intervalo  $\mu - \sigma$  a  $\mu + \sigma$ .
- Al menos  $3/4$  de las mediciones están en el intervalo  $\mu - 2\sigma$  a  $\mu + 2\sigma$ .
- Al menos  $8/9$  de las mediciones están en el intervalo  $\mu - 3\sigma$  a  $\mu + 3\sigma$ .



## Teorema de Chebyshev

---

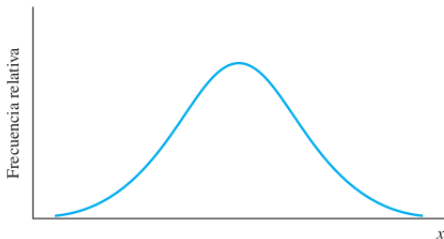
Ejemplo. La media y varianza de una muestra de  $n = 25$  mediciones son 75 y 100, respectivamente. Use el teorema de Chebyshev para describir la distribución de mediciones.

Solución Nos dan  $\bar{x} = 75$  y  $s^2 = 100$ . La desviación estándar es  $s = \sqrt{100} = 10$ . La distribución de mediciones está centrada alrededor de  $x = 75$ , y el teorema de Chebyshev establece que:

- Al menos  $3/4$  de las 25 mediciones están en el intervalo  $\bar{x} \pm 2s = 75 \pm 2(10)$ , esto es, 55 a 95.
- Al menos  $8/9$  de las mediciones están en el intervalo  $\bar{x} \pm 3s = 75 \pm 3(10)$ , esto es, 45 a 105.

## Regla empírica

Como el teorema de Chebyshev se aplica a cualquier distribución, es muy conservador. Ésta es la razón por la que hacemos hincapié en “al menos  $1 - (1/k^2)$ ” en este teorema. Otra regla para describir la variabilidad de un conjunto de datos no funciona para todos los conjuntos de datos, pero funciona muy bien para datos que “se apilan” en la conocida forma de montículo de la figura



# Regla empírica

---

---

**Regla empírica** Dada una distribución de mediciones que tiene forma aproximada de montículo:

El intervalo  $(\mu \pm \sigma)$  contiene aproximadamente 68% de las mediciones.

El intervalo  $(\mu \pm 2\sigma)$  contiene aproximadamente 95% de las mediciones.

El intervalo  $(\mu \pm 3\sigma)$  contiene aproximadamente 99.7% de las mediciones.

---

# Regla empírica

---

## Ejercicio

En un estudio de tiempo efectuado en una planta manufacturera, el tiempo para completar una operación especificada se mide para cada uno de los  $n = 40$  trabajadores. Se encuentra que la media y la desviación estándar son 12.8 y 1.7, respectivamente. Describa los datos muestrales usando la Regla empírica.

## Regla empírica

---

**Solución** Para describir los datos, calcule estos intervalos:

$$(\bar{x} \pm s) = 12.8 \pm 1.7 \quad \text{o} \quad 11.1 \text{ a } 14.5$$

$$(\bar{x} \pm 2s) = 12.8 \pm 2(1.7) \quad \text{o} \quad 9.4 \text{ a } 16.2$$

$$(\bar{x} \pm 3s) = 12.8 \pm 3(1.7) \quad \text{o} \quad 7.7 \text{ a } 17.9$$

De acuerdo con la Regla empírica, se espera que aproximadamente 68 % de las mediciones caigan en el intervalo de 11.1 a 14.5, aproximadamente 95 % caigan en el intervalo de 9.4 a 16.2, y aproximadamente 99.7 % caigan en el intervalo de 7.7 a 17.9.

## Ejercicio

---

Los maestros-estudiantes son capacitados para desarrollar planes de lecciones, en la suposición de que el plan escrito les ayudará a trabajar de manera satisfactoria en el salón de clases. En un estudio para evaluar la relación entre planes de lección escritos y su implementación en el salón de clases, se calificaron 25 planes de lección en una escala de 0 a 34 de acuerdo a una Lista de verificación de Plan de lección. Las 25 calificaciones se muestran en la tabla. Use el teorema de Chebyshev y la Regla empírica (si es aplicable) para describir la distribución de estas calificaciones de evaluación.

### Calificaciones para evaluación de Plan de lección

26.1	26.0	14.5	29.3	19.7
22.1	21.2	26.6	31.9	25.0
15.9	20.8	20.2	17.8	13.3
25.6	26.5	15.7	22.1	13.8
29.0	21.3	23.5	22.1	10.2

## Mediciones de posición relativa

---

A veces es necesario conocer la posición de una observación respecto a otras de un conjunto de datos.

### Definición

El **puntaje z muestral** es una medida de posición relativa definida por

$$\text{puntaje } z = \frac{x - \bar{x}}{s}$$

Un puntaje z mide la distancia entre una observación y la media, medidas en unidades de desviación estándar.



## Mediciones de posición relativa

---

Por ejemplo, suponga que la media y desviación estándar de los puntajes de examen (basados en un total de 35 puntos) son 25 y 4, respectivamente. El puntaje  $z$  para una calificación de 30 se calcula como sigue:

$$\text{puntaje } z = \frac{x - \bar{x}}{s} = \frac{30 - 25}{4} = 1.25$$

El puntaje de 30 está a 1.25 desviaciones estándar arriba de la media ( $30 = \bar{x} + 1.25s$ ).

# Mediciones de posición relativa

---

De acuerdo con el teorema de Chebyshev y la Regla empírica,

- al menos 75% y más probablemente 95% de las observaciones están a no más de dos desviaciones estándar de su media: sus puntajes  $z$  están entre  $-2$  y  $+2$ . *Las observaciones con puntajes  $z$  mayores a 2 en valor absoluto se presentan menos del 5% del tiempo y son consideradas un tanto improbables.*
- al menos 89% y más probablemente 99.7% de las observaciones están a no más de tres desviaciones estándar de su media: sus puntajes  $z$  están entre  $-3$  y  $+3$ . *Las observaciones con puntajes  $z$  mayores a 3 en valor absoluto se presentan menos del 1% del tiempo y son consideradas muy poco probables.*

# Percentil

---

Un percentil es otra medida de posición relativa y se usa con más frecuencia para conjuntos grandes de datos. (Los percentiles no son muy útiles para conjuntos pequeños de datos).

## Definición

Un conjunto de  $n$  mediciones de la variable  $x$  se ha reacomodado en orden de magnitud. El  $p$ -ésimo percentil es el valor de  $x$  que es mayor a  $p\%$  de las mediciones y es menor que el restante  $(100 - p)\%$ .

# Percentil

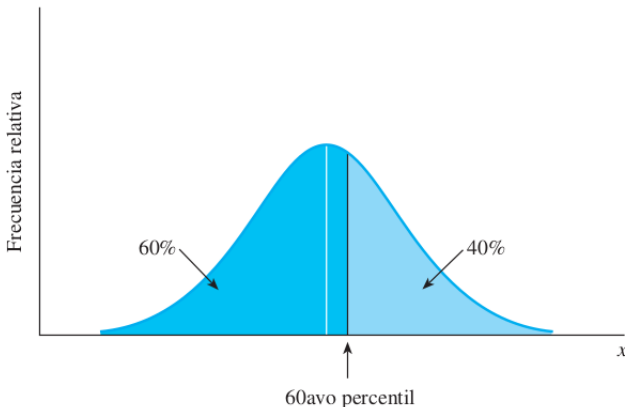
---

Supongamos que usted ha sido notificado que su calificación de 610, en el Examen verbal de graduación, lo ha colocado en el 60avo percentil en la distribución de calificaciones. ¿Dónde está su calificación de 610 en relación a las calificaciones de los otros que tomaron el examen?

Solución Calificar en el 60avo percentil significa que 60 % de todas las calificaciones de examen fueron más bajas que la calificación de usted y 40 % fueron más altas.

# Percentil

En general, el 60avo percentil para la variable  $x$  es un punto en el eje horizontal de la distribución de datos que es mayor a 60 % de las mediciones y menor que las otras.

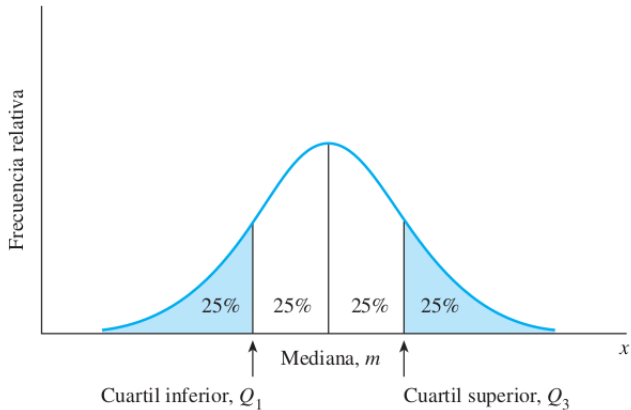


# Percentil-Cuartil

---

- La mediana es igual que el 50avo percentil.
- Los percentiles 25avo y 75avo, llamados cuartiles inferior y superior
- Veinticinco por ciento de las mediciones serán menores que el cuartil inferior (primero).
- 50 % serán menores que la mediana (el segundo cuartil).
- 75 % serán menores que el cuartil superior (tercero).

# Cuartil



# Cuartil

---

## Definición

Un conjunto de  $n$  mediciones en la variable  $x$  se ha acomodado en orden de magnitud. El **cuartil inferior (primer cuartil)**,  $Q_1$ , es el valor de  $x$  que es mayor a un cuarto de las mediciones y es menor que los restantes tres cuartos. El **segundo cuartil** es la mediana. El **cuartil superior (tercer cuartil)**,  $Q_3$ , es el valor de  $x$  que es mayor a tres cuartos de las mediciones y es menor que el restante un cuarto.



# Cuartil

---

## CÁLCULO DE CUARTILES MUESTRALES

- Cuando las mediciones están dispuestas en orden de magnitud, el **cuartil inferior**,  $Q_1$ , es el valor de  $x$  en la posición  $.25(n + 1)$ , y el **cuartil superior**,  $Q_3$ , es el valor de  $x$  en la posición  $.75(n + 1)$ .
- Cuando  $.25(n + 1)$  y  $.75(n + 1)$  no son enteros, los cuartiles se encuentran por interpolación, usando los valores de las dos posiciones adyacentes.<sup>†</sup>

# Cuartil

---

Encuentre los cuartiles inferior y superior para este conjunto de mediciones:

16, 25, 4, 18, 11, 13, 20, 8, 11, 9

**Solución** Ordene las  $n = 10$  mediciones de menor a mayor:

4, 8, 9, 11, 11, 13, 16, 18, 20, 25

Calcule

$$\text{Posición de } Q_1 = .25(n + 1) = .25(10 + 1) = 2.75$$

$$\text{Posición de } Q_3 = .75(n + 1) = .75(10 + 1) = 8.25$$

Como estas posiciones no son enteros, el cuartil inferior se toma como el valor  $3/4$  de la distancia entre la segunda y tercera mediciones ordenadas, y el cuartil superior se toma como el valor  $1/4$  de la distancia entre la octava y novena mediciones ordenadas. Por tanto,

$$Q_1 = 8 + .75(9 - 8) = 8 + .75 = 8.75$$

y

$$Q_3 = 18 + .25(20 - 18) = 18 + .5 = 18.5$$



# Intercuartil

---

Como la mediana y los cuartiles dividen la distribución de datos en cuatro partes, cada una de ellas conteniendo alrededor de 25 % de las mediciones,  $Q_1$  y  $Q_3$  son las fronteras superior e inferior para el 50 % central de la distribución. Podemos medir el rango de este “50 % central” de la distribución usando una medida numérica llamada rango intercuartil.

## Definición

El rango intercuartil (IQR) para un conjunto de mediciones es la diferencia entre los cuartiles superior e inferior; esto es,

$$IQR = Q_3 - Q_1 .$$

# El resumen de cinco números y la gráfica de caja

El **resumen de cinco números** consta del número más pequeño, el cuartil inferior, le mediana, el cuartil superior, y el número más grande, presentados en orden de menor a mayor:

**Min    $Q_1$    Mediana    $Q_3$    Max**

Por definición, un cuarto de las mediciones del conjunto de datos se encuentre entre cada uno de los cuatro pares adyacentes de números.

# El resumen de cinco números y la gráfica de caja

## PARA CONSTRUIR UNA GRÁFICA DE CAJA

- Calcule la mediana, los cuartiles superior e inferior y el IQR para el conjunto de datos.
- Trace una recta horizontal que represente la escala de medición. Forme una caja un poco arriba de la recta horizontal con los extremos derecho e izquierdo en  $Q_1$  y  $Q_3$ . Trace una recta vertical que pase por la caja en la ubicación de la mediana.

Una gráfica de caja se muestra en la figura 2.17.

