

Descriptive statistics Estadística descriptiva

Mario Enrique Rendón-Macías,¹ Miguel Ángel Villasís-Keever,¹ María Guadalupe Miranda-Novales²

Abstract

Descriptive statistics is the branch of statistics that gives recommendations on how to summarize clearly and simply research data in tables, figures, charts, or graphs. Before performing a descriptive analysis it is paramount to summarize its goal or goals, and to identify the measurement scales of the different variables recorded in the study. Tables or charts aim to provide timely information on the results of an investigation. The graphs show trends and can be histograms, pie charts, “box and whiskers” plots, line graphs, or scatter plots. Images serve as examples to reinforce concepts or facts. The choice of a chart, graph, or image must be based on the study objectives. Usually it is not recommended to use more than seven in an article, also depending on its length.

Keywords: Descriptive statistics; Tables; Figures; Graphics

Este artículo debe citarse como: Rendón-Macías ME, Villasís-Keever MÁ, Miranda-Novales MG. Estadística descriptiva. Rev Alerg Mex. 2016;63(4):397-407

¹Instituto Mexicano del Seguro Social, Hospital de Pediatría, Centro Médico Nacional Siglo XXI, Unidad de Investigación en Epidemiología Clínica. Ciudad de México, México

²Instituto Mexicano del Seguro Social, Hospital de Pediatría, Centro Médico Nacional Siglo XXI, Unidad de Investigación en Epidemiología Hospitalaria. Ciudad de México, México

Correspondencia: Mario Enrique Rendón-Macías.
drmariorendon@gmail.com

Recibido: 2016-09-29
Aceptado: 2016-10-02



Resumen

La estadística descriptiva es la rama de la estadística que formula recomendaciones de cómo resumir, de forma clara y sencilla, los datos de una investigación en cuadros, tablas, figuras o gráficos. Antes de realizar un análisis descriptivo es primordial retomar el o los objetivos de la investigación, así como identificar las escalas de medición de las distintas variables que fueron registradas en el estudio. El objetivo de las tablas o cuadros es proporcionar información puntual de los resultados. Las gráficas muestran las tendencias y pueden ser histogramas, representaciones en "pastel", "cajas con bigotes", gráficos de líneas o de puntos de dispersión. Las imágenes sirven para dar ejemplos de conceptos o reforzar hechos. La selección de un cuadro, gráfico o imagen debe basarse en los objetivos del estudio. Por lo general no se recomienda usar más de siete en un artículo destinado a una publicación periódica, parámetro que está también en función de la extensión misma del artículo.

Palabras clave: Estadística descriptiva; Cuadros; Figuras; Gráficas

Estadística descriptiva

El objetivo final de cualquier investigación es proporcionar evidencia objetiva suficiente para apoyar o refutar la o las hipótesis planteadas.^{1,2} La evidencia obtenida mediante la recolección planeada y cuidadosa de una investigación tiene que traducirse en datos o cifras. Al integrar y dar coherencia a los resultados de un trabajo, el investigador debe tener la capacidad de resumir y presentar datos de manera ordenada, sencilla y clara, para que puedan ser interpretados tanto por otros investigadores como por los revisores y lectores.

La estadística descriptiva es la rama de la estadística que formula recomendaciones sobre cómo resumir la información en cuadros o tablas, gráficas o figuras.^{2,3}

Prerrequisitos

Antes de realizar un análisis descriptivo sobre los resultados de una investigación, es fundamental revisar el o los objetivos de la misma. No es infrecuente que para ese momento se olviden los propósitos que llevaron a la realización del estudio, lo cual puede derivar en un trabajo infructuoso y con alto riesgo de generar confusión más que conclusiones acertadas.⁴

Una vez recuperados el o los objetivos de la investigación es necesario identificar las escalas de medición de las variables registradas, las cuales para fines prácticos pueden ser:

- Cuantitativas (métricas).
- Cualitativas (categóricas).

Las primeras se definen por la existencia de una unidad de medición, que puede ser contable (unidades enteras), medible o ponderada por algún atributo físico con algún instrumento. Asimismo, pueden ser clasificadas como continuas si aceptan fracciones, o discretas si solo consideran unidades enteras. Ejemplos de estos tipos de variables son el volumen de espiración medido en litros (variable cuantitativa continua) y el número de respiraciones por minuto (variable cuantitativa discreta).²

Las variables cualitativas se caracterizan por clasificar a los individuos o fenómenos solo con relación a sus atributos. Pueden ser nominales cuando los atributos usados son únicos para una condición (excluyentes mutuamente) y solo existen posibilidades conocidas (exhaustivas). Si el atributo solo acepta dos condiciones, las variables reciben el nombre de *nominales dicotómicas*, pero si hay más posibilidades se les denomina *nominales politómicas*.^{2,3} Un ejemplo de variable dicotómica sería la presencia o ausencia de eccema, mientras que de una variable politómica, los tipos de alergia.

Por otro lado, cuando la clasificación cualitativa se basa en un orden jerárquico de los atributos, a las variables se les conoce como *ordinales*, tal como sucede con la dificultad respiratoria, la cual puede ser ausente, leve, moderada o severa.⁵

Medidas de tendencia central y de dispersión

En general, para resumir o presentar los datos obtenidos de un proyecto de investigación inicialmente se debe tratar de ubicar cómo se distribuyen, lo cual se realiza de acuerdo con la escala de medición de cada variable.

Escala cuantitativa

Algunos datos deben resumirse en un estimador de promedio y otros en uno de dispersión. El estimador de promedio indica la tendencia central o cifra que representa mejor el valor de la muestra, la cual puede ser:

- Promedio o media (aritmética), obtenido con la suma de todos los valores individuales entre el número total de valores; representa el punto de equilibrio de la distribución de los datos.
- Mediana, que representa la cifra o valor que divide la muestra en dos mitades, es decir, el valor donde 50% de la población está por debajo o arriba del mismo.
- Moda o valor más frecuentemente encontrado en las mediciones.²⁻⁴

Las medidas de dispersión para las variables cuantitativas son tres: la desviación estándar o desviación típica, los rangos intercuartílicos y los valores mínimo y máximo. Todas permiten entender cómo se alejan los datos del promedio y la distribución dentro de los límites medidos. No se abordan las fórmulas y cálculos que se realizan para obtenerlas dado que se describen en diferentes libros de texto e incluso existen programas de computación.²⁻⁴

Sin embargo, antes de proceder a calcular cualquiera de las medidas descritas, es necesario que primero se establezca la distribución de los datos que se están analizando. Para fines prácticos, para los datos cuantitativos se debe establecer si tienen una distribución normal o gaussiana (es decir, que se asemeje a la curva de distribución normal o curva de Gauss). La prueba Kolmogorov-Smirnov es una de las pruebas estadísticas para determinar si la distribución de un grupo de datos es normal, aunque también se puede utilizar el sesgo y la curtosis, o bien, gráficos tipo p-p o q-q.

Si se determina que la distribución de los datos es normal, entonces los resultados solamente se

pueden indicar con dos estimadores: media y desviación estándar. En una distribución normal, 95% de los datos estudiados se encuentra dentro de ± 2 desviaciones estándar a partir de la media.

Cuando la distribución no es normal, resulta más informativo mostrar los valores percentílicos o los cuartílicos. Los valores percentílicos establecen la probabilidad de 100% de encontrar un valor, los cuartílicos dividen el total de los datos en cuatro porciones equivalentes a 25% de los datos. En una división percentilar, el percentil 5 indica la cifra donde 5% de los datos está por debajo de esta y 95% arriba de la misma. Si se anotan los percentiles 5 y 95, se informa el rango de valores en el que se encuentra 90% de la distribución.

En los cuartiles suele informarse el Q1 (cuartil 1) y el Q3 (cuartil 3); el primero equivale al percentil 25 y el segundo, al percentil 75. Los valores que se encuentran en el intervalo de Q1 y Q3 dan cuenta de 50% de los datos de la distribución más cercanos a la mediana.

La última medida de dispersión es el rango o intervalo entre el valor mínimo y máximo, el cual se obtiene de la resta del valor mayor menos el valor menor + 1.^{1,2}

Escala cualitativa

Los datos deben ser resumidos en frecuencias simples o relativas en porcentaje. La frecuencia simple solo es el conteo de los eventos en cada categoría. La frecuencia relativa se obtiene dividiendo cada conteo de eventos de esa categoría entre el total de las mediciones. Se pueden presentar los resultados en fracciones o multiplicarse por 100, para expresarlos en porcentajes. Esta última opción suele ser más fácil para su interpretación.

Cuadros o tablas

Estos consisten en matrices de datos que permiten determinar cifras puntuales sobre las mediciones realizadas. Constan de tres partes fundamentales: el título, el cuerpo (cabecera de tabla y matriz de datos) y los acotamientos o aclaraciones.²

El título debe informar sobre el contenido del cuadro, las variables mostradas y el número de sujetos o unidades de estudio.

La primera fila del cuadro mostrará los encabezados de cada columna. Estos deben informar sobre los valores estadísticos usados para resumir

los datos de cada una de las variables (porcentajes, promedio, desviación estándar, etcétera). En las filas de la primera columna se anotan las variables consideradas en el cuadro (Figura 1).

Cuando un estudio pretende comparar dos o más grupos, en las columnas se anotarán estos; la recomendación es indicar primero el grupo experimental o de interés y después el o los controles.

Es importante ordenar las variables en un sentido lógico o secuencial, es decir, en un cuadro que descri-

ba datos clínicos se deberá anotar primero los datos de los síntomas y signos, seguidos de los datos de laboratorio, imagen, tratamientos, resultados, etcétera.

Por último, al incorporar los datos en cada columna se recomienda que tengan una secuencia relacionada, generalmente con la frecuencia del evento presentado, es decir, iniciando con el valor mayor (100%, 60%, 20%) o con el valor menor (5, 8, 10, 15, 20). Otra forma de presentar los datos de manera ordenada puede ser de acuerdo con las categorías de

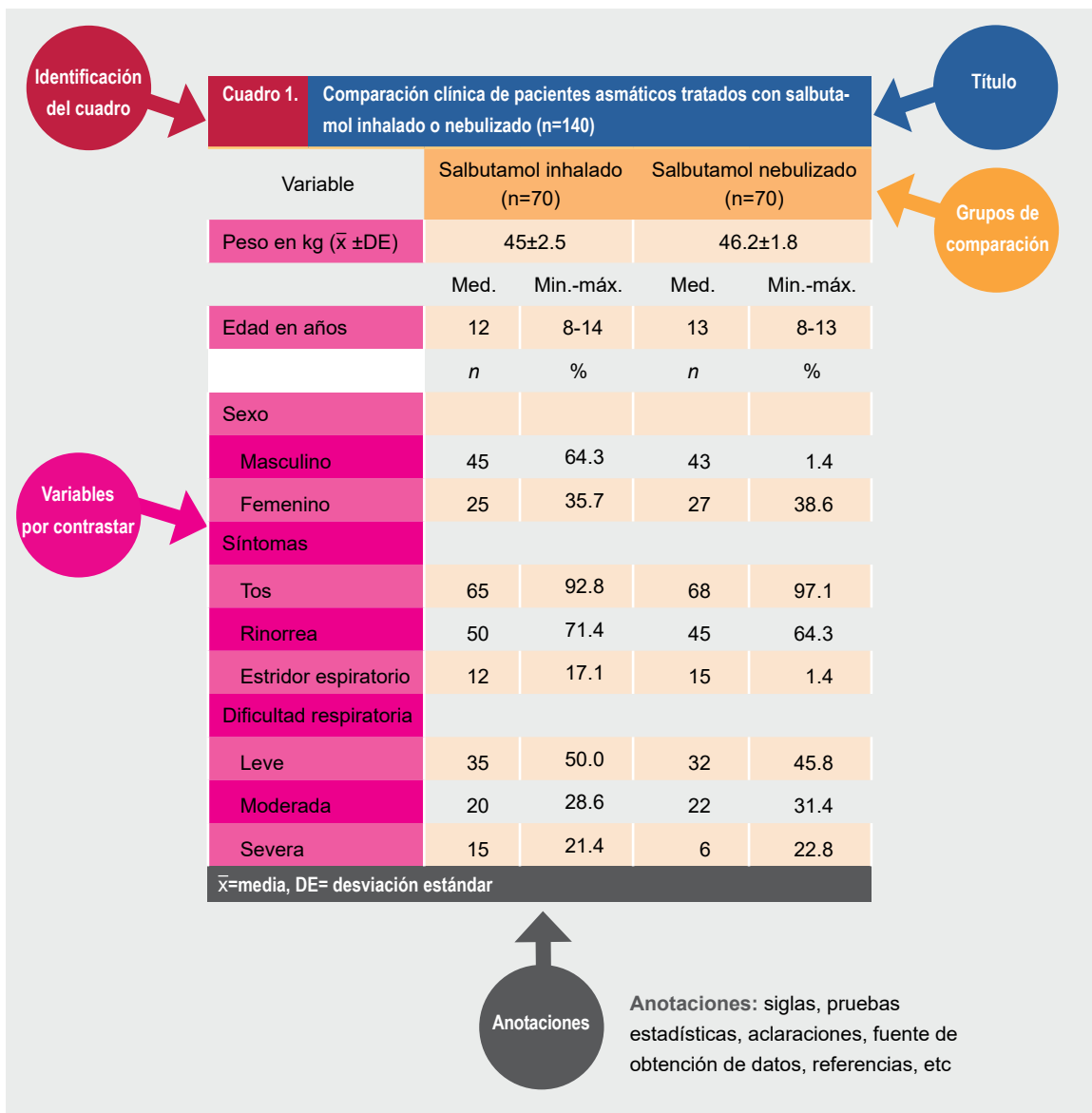


Figura 1. Esquema de un cuadro o tabla.

cada variable que se señala en la primera columna (por ejemplo, estadio I, estadio II, estadio III), aun cuando los datos de las categorías no sigan un orden progresivo.

Gráficas

Las gráficas tienen como objetivo mostrar tendencias más que datos puntuales. También son muy útiles para comparar visualmente los resultados de los grupos; sobre todo se emplean para resaltar hallazgos o resultados importantes. Aunque en la actualidad existen programas para realizar imágenes complejas, es preferible emplear gráficos de relación entre dos o tres variables. Una gráfica para la distribución de una sola variable solo es útil cuando esta es la respuesta principal de una investigación, de lo contrario resulta superflua y la información bien puede ser expresada textualmente en los resultados de un artículo.

Cuando se construye una gráfica hay tres aspectos fundamentales por considerar:

- La identificación clara de las variables.
- La descripción de la o las escalas utilizadas.
- El uso de la menor cantidad posible de palabras, pero suficientes para facilitar la comprensión.^{1,3-6}

Con excepción de las gráficas de “pastel”, las restantes suelen tener al menos dos ejes de construcción: el eje de la ordenada o X, que suele ser la variable independiente o predictora y el eje de la abscisa o Y, que corresponde a la dependiente o el resultado. La decisión de cuál elegir dependerá del número de variables por mostrar y de la escala de medición de cada una.

Gráficas para una sola variable

Cuando se desea señalar un resultado importante en la distribución de una sola variable existen al menos dos opciones muy útiles. Si la variable es medida en una escala cuantitativa es común usar los histogramas (Figura 2). Esta gráfica puede ser construida a partir de una variable con escala de medición discreta o continua. Se caracteriza por la unión entre las barras asumiendo que existe un cambio en las frecuencias al pasar un umbral, el cual es claramente definido en medidas de conteo, pero es arbitrario (por lo general al valor 0.5) en las escalas continuas. En ocasiones se pueden unir los puntos centrales del intervalo de una clase para formar un polinomio.

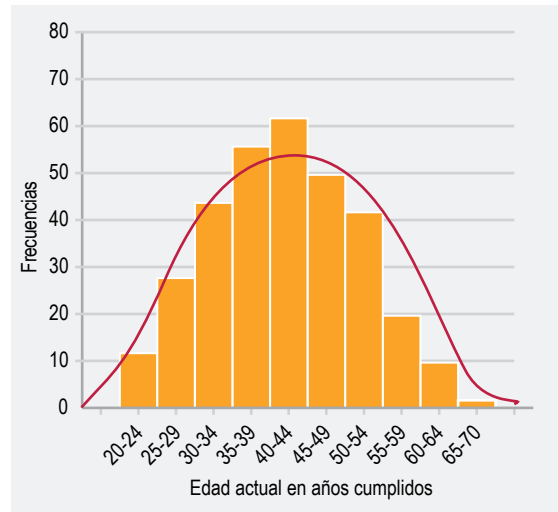


Figura 2. Ejemplo de un histograma. La línea representa la probabilidad de mostrar una distribución normal o en forma de campana.

En los resultados obtenidos de la práctica clínica se utilizan poco estas gráficas.^{4,6,7}

En la investigación clínica es más común resumir una sola variable cualitativa en una gráfica de “pastel” (Figura 3). En esta modalidad cada porción se interpreta como “la proporción que ocupa por su frecuencia dentro de un todo”. Si todos los pacientes la tienen, entonces equivale a toda el área del pastel, pero cuando hay más de un atributo, la proporción se representa como el tamaño de una “rebanada” de este. Aunque hay fórmulas para calcular los grados equivalentes para cada proporción,² con diversos programas de computadora (por ejemplo, Excel) se obtienen en forma automática. En estas gráficas de “pastel” se recomienda que existan más de dos categorías y no incorporar porciones menores a 1% o en su defecto agruparlas con otras también pequeñas. Por otro lado, es conveniente ordenar las rebanadas de mayor a menor en el sentido de las manecillas de un reloj, a menos que exista un orden preestablecido por gravedad o por cuestiones clínicas.

Gráficas para dos o más variables no relacionadas

Con ellas se pretende comparar resultados entre grupos. En un intento por hacer la comparación sencilla y clara, si las variables son cualitativas se recomienda usar gráficas de barras (Figura 4), en las cuales es

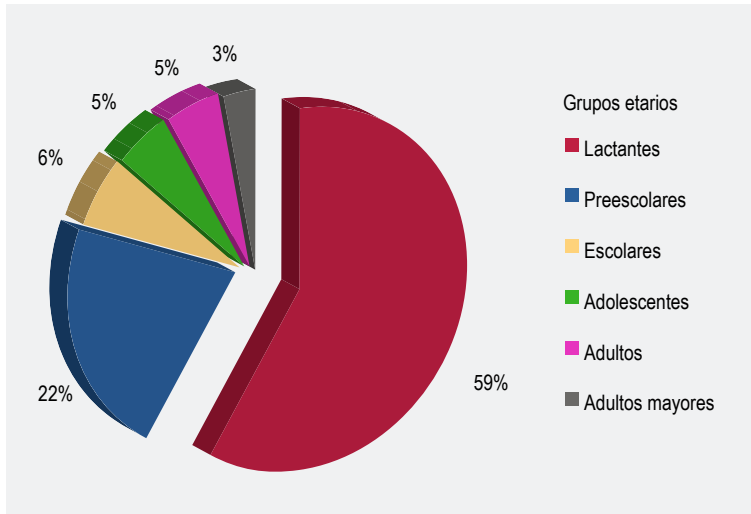


Figura 3. Ejemplo de gráfica tipo "pastel". Distribución por grupos etarios de alergía a la proteína de la leche (N= 2245).

aconsejable colocar las variables cualitativas en el eje de las X y la frecuencia de presentación de cada categoría en el eje de las Y, ya sea en su valor real o relativo (proporción o porcentaje). Es importante resaltar que se debe graficar los valores de frecuencias reales solo cuando los grupos de comparación tienen el mismo número de pacientes, de lo contrario se deben representar las proporciones o los porcentajes.^{1,6,7} Como se muestra en las Figuras 4a y 4b, la

comparación de frecuencias simples en grupos de tamaño diferente puede distorsionar la realidad.

Cuando una variable es cualitativa (por ejemplo, un grupo con rinitis alérgica contra un grupo control) y la variable de desenlace es cuantitativa (por ejemplo, el número de eosinófilos por mm³ de suero), los resultados pueden representarse con distintos gráficos (Figuras 5 y 6). La decisión depende de dos factores principales: la distribución de los da-

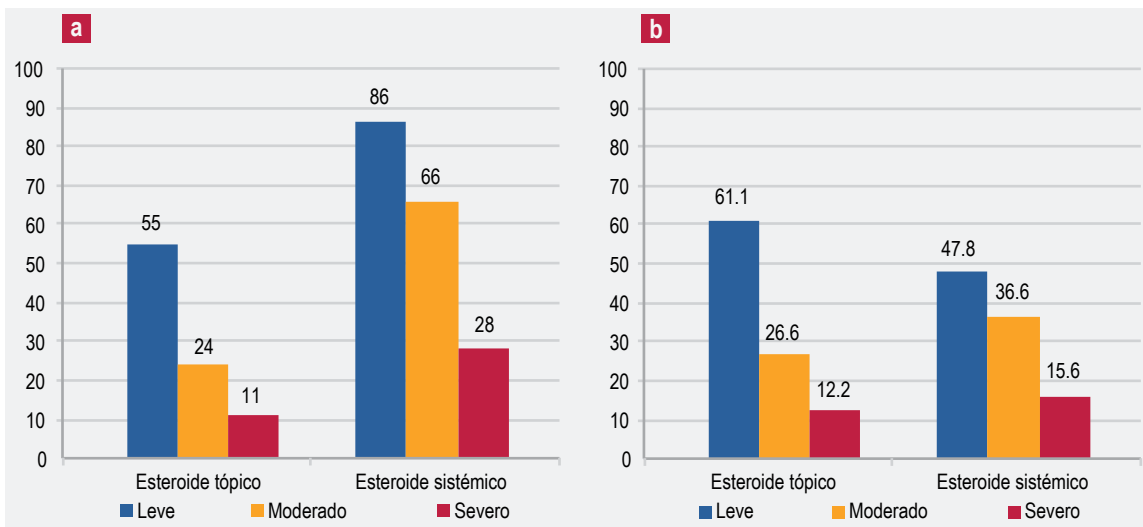


Figura 4. Ejemplo de gráfica tipo "barras" en grupos comparativos. Distribución por severidad de los síntomas entre pacientes con tratamiento tópico (n=90) y sistémico con esteroides (n=180).

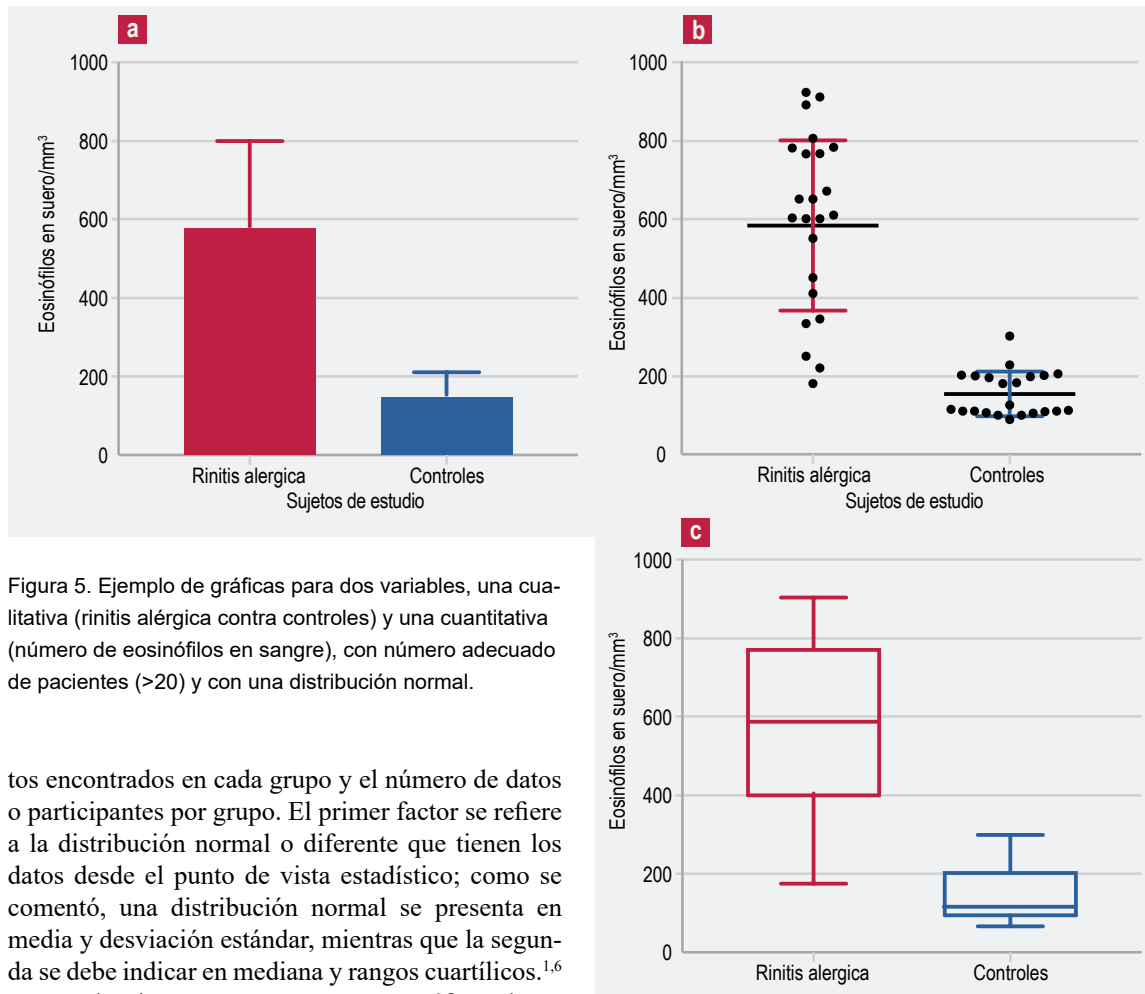


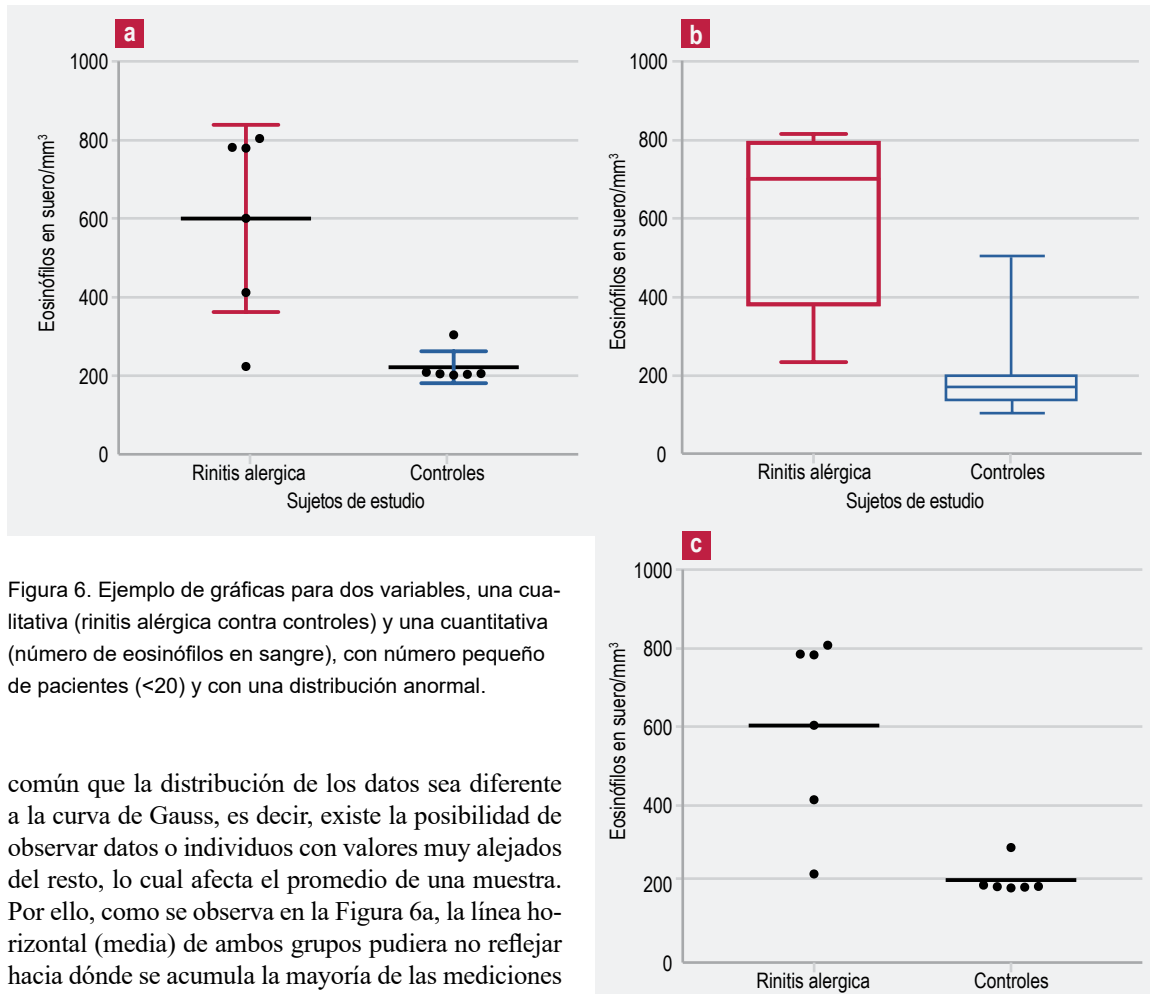
Figura 5. Ejemplo de gráficas para dos variables, una cualitativa (rinitis alérgica contra controles) y una cuantitativa (número de eosinófilos en sangre), con número adecuado de pacientes (>20) y con una distribución normal.

tos encontrados en cada grupo y el número de datos o participantes por grupo. El primer factor se refiere a la distribución normal o diferente que tienen los datos desde el punto de vista estadístico; como se comentó, una distribución normal se presenta en media y desviación estándar, mientras que la segunda se debe indicar en mediana y rangos cuartílicos.^{1,6}

En la Figura 5 se muestran tres gráficas de un mismo resultado. Si la distribución es normal puede usarse cualquiera (como en el caso que se muestra). La Figura 5a solo señala una columna que termina en la media de cada grupo y la barra con la línea horizontal equivale a la distancia de una desviación estándar. Esta solo indica una desviación estándar por arriba del promedio, pero se espera que sea la misma hacia abajo. La Figura 5b muestra el valor de cada paciente según su tratamiento y el valor promedio con una línea horizontal negra; con barras horizontales más pequeñas se marca una desviación estándar arriba y debajo de la media. Como se observa, esta gráfica se ve saturada y no aporta nada distinto a la anterior. La Figura 5c es una gráfica de cajas, la cual muestra una raya interior horizontal que corresponde al valor de la mediana, el borde inferior de la caja es el percentil 25 o cuartil 1 (Q1) y el borde

superior equivale al percentil 75 o cuartil 3 (Q3). Los bigotes o líneas horizontales superior e inferior equivalen a los percentiles 10 y 90. Esta figura se ve muy simétrica a partir de la mediana, lo cual traduce una distribución muy cercana a la normal, por ello no aporta nada nuevo a la Figura 5a. Esta gráfica es muy útil para la comparación de grupos cuando la distribución no es normal, es decir, cuando no hay simetría de las cajas respecto a la mediana.

Como señalamos, el segundo factor que influye para decidir qué gráfica debe elaborarse es el número de pacientes estudiados. En el ejemplo anterior, este número era suficiente para que los datos alcanzaran una distribución normal. En la Figura 6 mostramos qué sucede cuando el número de datos o pacientes es pequeño (<20). En estas condiciones es



Gráficas para dos variables cuando hay seguimiento

En ocasiones encontramos fenómenos que cambian con el tiempo o de los cuales se desea saber cómo influye este.⁸ En las gráficas para dos variables cuando hay seguimiento, la variable de tiempo debe ser

graficada en el eje de las X porque la lectura es de izquierda a derecha, además de que con ello se da la sensación de progresión. En el eje de las Y se anota la variable de desenlace en la escala medida (Figura 7). Aunque, el solo uso de barras puede ser suficiente para apreciar un recorrido en el tiempo (Figura 7a), es recomendable dibujar una línea que una los momentos (promedios) entre los eventos. Esta línea permite entender mejor la relación entre los diferentes momentos evaluados (Figura 7b).

Como se comentó previamente, las gráficas anteriores pueden ser modificables según la distribución de los datos y el número de pacientes o medidas realizadas.²⁻⁴ En particular cuando se trata de pocos pacientes se puede construir una gráfica de unión antes y después de cada evaluación individual (Fi-

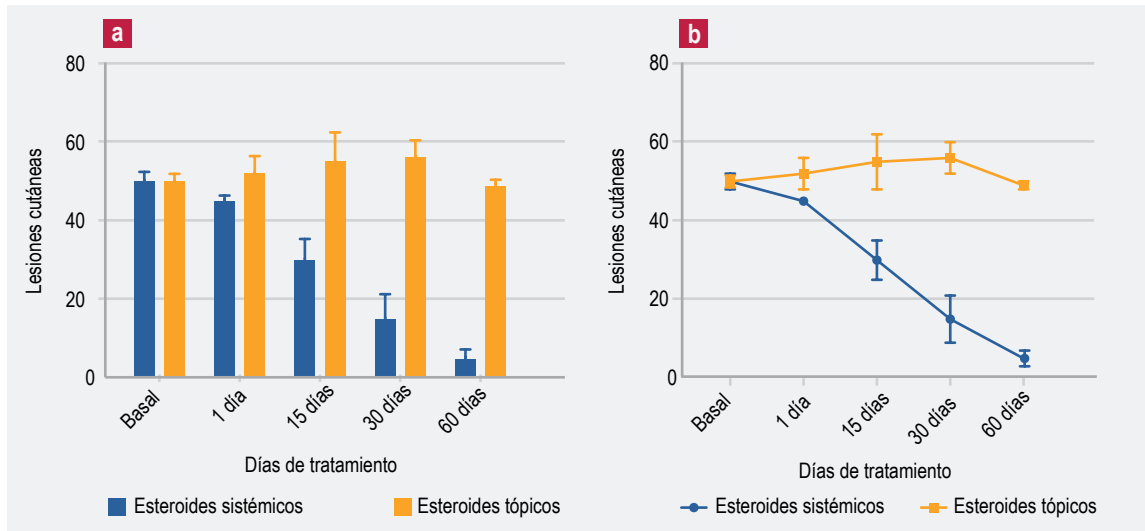


Figura 7. Ejemplo de gráfica de secuencia en el tiempo de una medición repetida. Las líneas en la figura b une las medias (bolos) y las líneas verticales la distancia de una desviación estándar por arriba y debajo de esta.

gura 8). Esta gráfica permite establecer el comportamiento individualizado para muestras pequeñas.

Gráficas entre dos variables cuantitativas (correlación)

Este tipo de gráficas se utiliza cuando la intención es correlacionar dos mediciones métricas.^{1,3} Cada sujeto tiene dos valores y este par es graficado en dos ejes (X y Y). Normalmente, el eje X representa la variable independiente cuando se trata de dar una dirección a la relación entre estas dos variables, y el eje de las Y representa la variable dependiente (Figura 9). En este tipo de gráficas es muy importante anotar las escalas y unidades de medición de cada variable. Se recomienda iniciar con 0 en el sitio de intersección de ambos ejes, sin embargo, se puede iniciar una escala en cualquier cifra, siempre y cuando se anote claramente la cifra de inicio. Es recomendable anotar: 0 -// - cifra de inicio. Así mismo, es recomendable que el tamaño de los ejes sea igual, para una mejor interpretación de la relación.⁴ La relación entre las variables es observable en la tendencia de los puntos hacia la formación de una nube o una línea. Entre más clara sea dibujada una línea la correlación será más adecuada.^{4,6}

Para facilitar la interpretación de esta gráfica por parte de los lectores se recomienda dibujar la línea de mejor ajuste de los puntos. Para encontrar esta lí-

nea existen diferentes métodos descritos en diversos libros especializados, aunque también están disponibles en distintos programas estadísticos.²

Por último, se pueden graficar en la misma figura la correlación de más de un grupo o categoría de la misma variable, un ejemplo podría ser la correlación entre número de lesiones y dosis en “mujeres” y en “hombres”. Lo recomendable es no saturar con

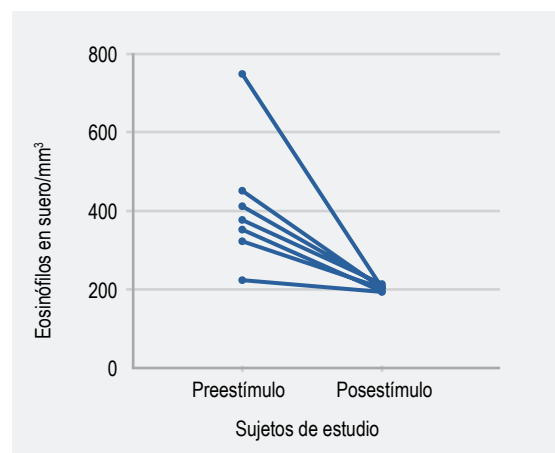


Figura 8. Ejemplo de una gráfica de comparación antes y después de un momento. Cada línea une un valor previo y posterior al mismo de un paciente individual.

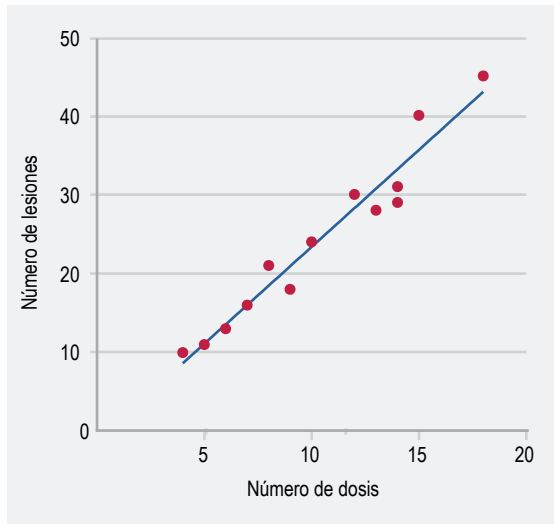


Figura 9. Ejemplo de una figura de puntos de correlación entre dos variables cuantitativas. La línea refleja la tendencia de la relación, en este caso una correlación positiva.

múltiples puntos y líneas para una interpretación más sencilla.

Gráficas de supervivencia

Este tipo de gráfica es muy común en la actualidad. En ella, el eje de la Y representa la probabilidad (0 a

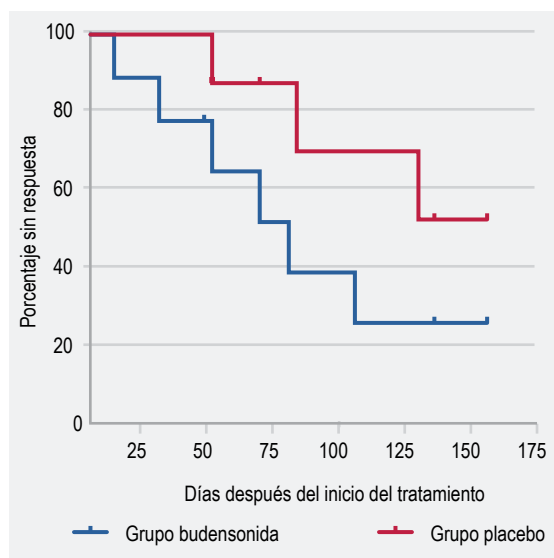


Figura 10. Ejemplo de curva de Kaplan y Meier o gráfica de supervivencia.

1 o en porcentaje de 0 a 100%) de continuar en una condición clínica con el tiempo. Cada vez que sucede un evento evaluado (por ejemplo, muerte, recaída, recurrencia, etcétera), la curva cae en proporción con los individuos que la padecen respecto al total seguido hasta ese momento. Cuando un paciente se pierde o sale del estudio por otro motivo distinto al evento de interés, se marca con una línea vertical sobre la curva sin modificar la probabilidad. El tipo de gráfica más utilizado es la curva de supervivencia de Kaplan-Meier (Figura 10).^{9,10}

Figuras o imágenes

Estas consisten en iconografías que pretenden mostrar resultados específicos. En particular son muy útiles cuando la descripción de un fenómeno con textos se hace compleja (“una imagen dice más que mil palabras”) o cuando se desea generar una imagen que recuerde el lector. Pueden ser esquemas o fotografías. Lo importante en una publicación es buscar la más representativa del fenómeno que se desea mostrar. Siempre es conveniente que sea de alta calidad y, como en los casos anteriores, no debe faltar un pie de figura que explique el propósito de la imagen y puntualice los detalles que deberán observarse. Se pueden utilizar flechas dentro de las imágenes para señalarlos. Si se usa un esquema, es importante que sea lo más sencillo posible, reduciendo al máximo las palabras y las flechas. Si utiliza colores, que estos sean contrastantes y de preferencia no más de tres.

Por último, en la actualidad hay muchas gráficas utilizadas en algunos estudios particulares como los metaanálisis.¹¹ Todos requieren conocer la metodología y estadística específica. En esta revisión solo presentamos la estadística descriptiva más utilizada en publicaciones de trabajos de investigación.

Conclusiones

La estadística descriptiva tiene como objetivo resumir la evidencia encontrada en una investigación de manera sencilla y clara para su interpretación. Consta de tablas o cuadros, figuras o gráficas e imágenes o fotografías. Los cuadros se utilizan para resumir datos y mostrar cifras puntuales. Las figuras o gráficas tienen la finalidad de señalar tendencias y comparaciones. Las imágenes o fotografías permiten mostrar fenómenos difícilmente explicables en el texto. Se recomienda no usar más siete de estas herramientas en una publicación.

Referencias

1. Spriestersbach A, Röhrig B, du Prel JB, Gerhold A, Blettner M. Descriptive statistics. *Dtsch Arztebl Int* 2009;106(36):578-583. doi: 10.3238/arztebl.2009.0578
2. Diggle PJ, Chetwynd AG. *Statistics and scientific method. An introduction for students and researchers.* UK: Oxford University Press; 2013. p. 36-56.
3. Sonnad S. Describing data: Statistical and graphical methods. *Radiology.* 2002;225(3):622-628.
4. Sperandei S. The pits and falls of graphical presentation. *Biochem Med (Zagreb).* 2014;24(3):311-320. doi: 10.11613/BM.2014.033.
5. Villasís-Keever MA, Miranda-Navales MG. El protocolo de investigación IV: las variables de estudio. *Rev Alerg Mex.* 2016;63:303-310.
6. Lang T. Twenty statistical errors even you can find in biomedical research articles. *Croat Med J.* 2004;45(4):361-370.
7. Overholser BR, Sowinski KM. Biostatistics primer: Part I. *Nut Clin Pract.* 2007;22(6):629-635.
8. Tobías A, Sáez M, Galán I. Herramientas gráficas para el análisis descriptivo de series temporales en la investigación médica. *Med Clin (Barc).* 2004;122(18):701-706.
9. George B, Seals S, Aban I. Survival analysis and regression models. *J Nucl Cardiol.* 2014;21(4):686-694. doi: 10.1007/s12350-014-9908-2
10. Rich JT, Gail-Neely J, Paniello RC, Voelker CCJ, Nussenbaum B, Wang EW. A practical guide to understanding Kaplan-Meier curves. *Otolaryngol Head Neck Surg.* 2010;143(3):331-336. doi: 10.1016/j.otohns.2010.05.007
11. Ressing M, Blettner M, Klug SJ. Systematic reviews and meta-analyses. *Dtsch Arztebl Int.* 2009;106(27):456-463. doi: 10.3238/arztebl.2009.0456