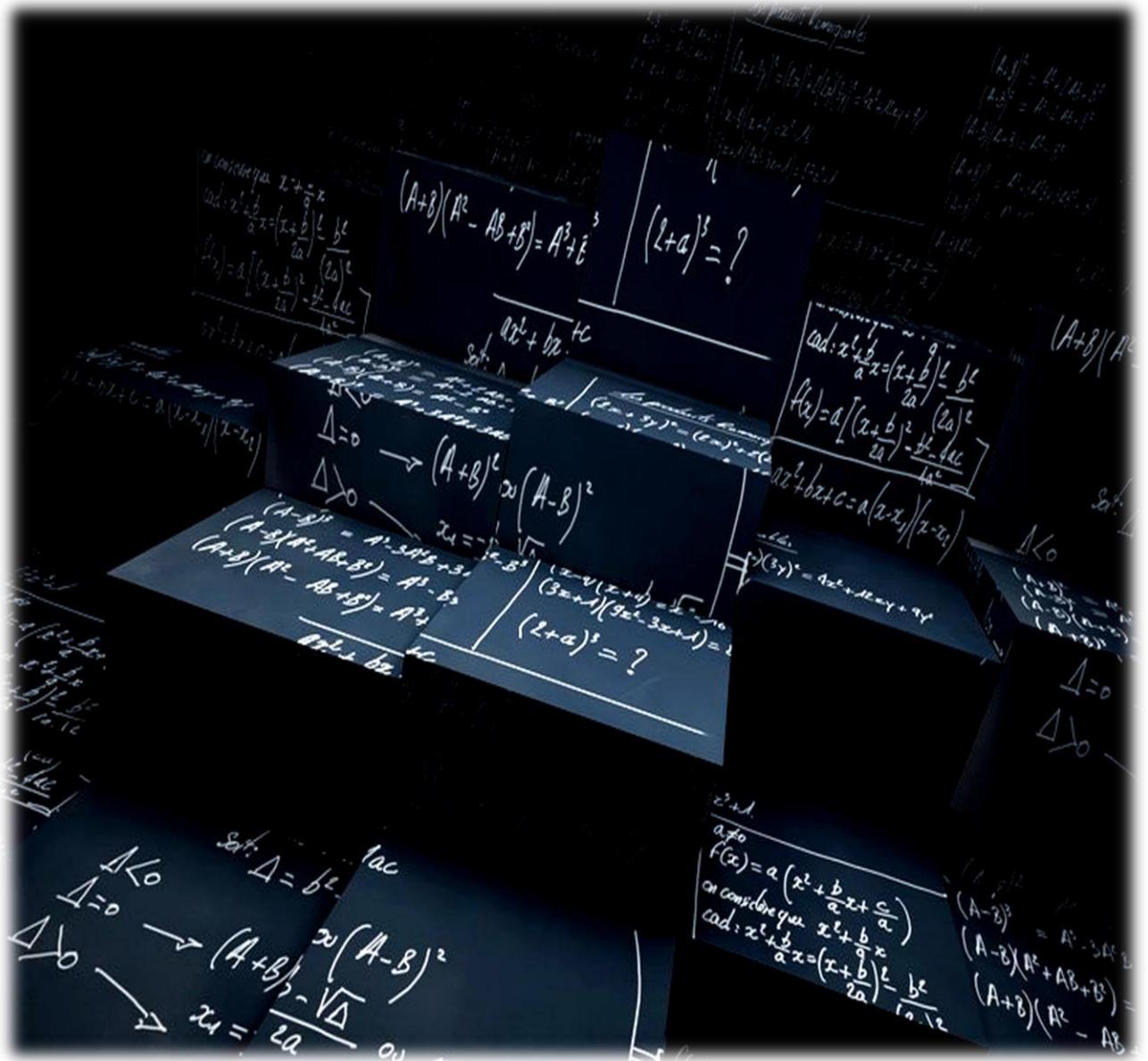


Tutoriales de Estadística



Estadística y Sistemas de Información Educativa
Escuela de Ciencias de la Educación
Facultad de Filosofía y Humanidades
Universidad Nacional de Córdoba

Tutorial de Estadística

Cátedra

**Estadística y Sistemas de Información Educativa
Escuela de Ciencias de la Educación
Facultad de Filosofía y Humanidades
Universidad Nacional de Córdoba**

Lic. Jorge Lorenzo

Este libro está bajo la siguiente licencia Creative Commons



Disponible en repositorio Ansenza FFyH. UNC.

<https://ansenuza.unc.edu.ar/>

Índice

Introducción: página 7

Capítulo 1

Variables, Unidad de Análisis: página 11

Sistemas de medición: página 14

Otros usos del término escala: página 21

Capítulo 2

Técnicas de muestreo: página 23

Muestra y Población: página 23

Población: página 23

Muestra e inferencia: página 25

Muestreo probabilístico: página 27

Muestreo aleatorio simple: página 28

Muestreo aleatorio sistemático: página 30

Muestreo aleatorio estratificado: página 32

Muestreo por conglomerados: página 36

Muestreo No Probabilístico: página 38

Muestreo accidental: página 38

Muestreo por cuotas: página 38

Muestreo intencional: página 39

Bola de nieve: página 39

Comentarios Finales: página 39

Misceláneas página 40

Capítulo 3

Tablas de frecuencia: página 43

Tabla de frecuencia: sistema de medición nominal: página 43

Utilidad de la tabla de frecuencia a través de un ejemplo: página 46

Tabla de frecuencia: sistema de medición ordinal: página 47

Tablas de frecuencias para variables métricas: página 51

Representación gráfica de las tablas de frecuencias: página 54

Diagrama de barras: página 54

Diagrama de sectores: página 55

Diagrama de barras agrupadas: página 55
Histograma: página 56
Polígonos de frecuencia: página 59
Gráfico de ojiva: página 61
Otros tipos de gráficas: página 62
Pirámides de población: página 62
Pictogramas: página 66
Mapas: página 67
Gráficas polares: página 68
Box Plot o Diagramas de cajas: página 69
Diagrama de cajas y medidas de posición: cuartiles: página 70
Diagrama de cajas – casos con valores atípicos o extremos: página 71
Diagrama de cajas – restricción de la variabilidad: página 73
Diagrama de cajas en distribuciones sesgadas: página 74

Capítulo 4

Medidas de tendencia central: página 77
Modo: página 77
Mediana: página 80
Media Aritmética o Promedio: página 81
Varianza y Desviación Estándar: página 82
Coeficiente de variación de Pearson: página 85
Procedimiento para el cálculo del promedio, la varianza y la desviación estándar: página 86

Capítulo 5

La distribución normal: página 88
La distribución normal como modelo matemático: la normal estandarizada: página 91
La distribución normal estandarizada y la proporción de casos: página 94
Áreas bajo la curva como distribución de probabilidades: página 96
Misceláneas: página 102
Distribuciones asimétricas: página 104
Medidas de asimetría y curtosis de una distribución: página 106

Capítulo 6

Estimación de parámetros: página 110
Estimación de una media poblacional: página 110
Estimación de una proporción poblacional: página 112
Comparación de proporciones: página 114
Comparación de proporciones muestrales: página 115
Estimación de parámetros: conceptos teóricos: página 117

Capítulo 7

- Prueba de hipótesis sobre la media paramétrica: página 118
- Los datos y el modelo estadístico: página 119
- La distribución normal como modelo estadístico: página 125
- Comparación de distintos promedios con el modelo estadístico de la distribución normal estándar: página 129
- Prueba de hipótesis sobre media paramétrica: página 130

Capítulo 8

- Prueba de hipótesis y el modelo estadístico χ^2 : página 137
- Cálculo del estadístico χ^2 a partir de la tabla de contingencia: página 138
- Reglas de decisión basadas en χ^2 y errores de decisión: página 141
- El grado de asociación entre las variables: página 143
- Coficiente de contingencia C de Pearson: página 143
- Coficiente V de Cramer: página 144
- Ejemplo de aplicación: página 144
- Las frecuencias esperadas como eventos independientes: página 146
- Los grados de libertad en tablas de contingencia: página 147
- Distribución χ^2 : página 149
- Recomendaciones para el uso de la prueba χ^2 : página 150
- La prueba de la mediana: página 150
- La mediana combinada: página 151
- Análisis de una muestra simple: prueba χ^2 para la bondad de ajuste: página 153

Capítulo 9

- Correlación de variables: página 156
- El diagrama de dispersión: página 157
- Coficiente de correlación r de Pearson: página 161
- Coficiente de correlación r_s de Spearman: página 165
- Coficiente de correlación y prueba de hipótesis: página 166
- Algunas precisiones sobre el coeficiente de correlación: página 171
- Otras medidas de asociación simétrica y asimétrica

Introducción

La estadística es una rama de las matemáticas y como tal, nos ofrece distintos modelos de análisis de datos cuantitativos. Puede ser vista como una herramienta para el resumen de grandes volúmenes de información y un puntal para la investigación científica. El lenguaje estadístico establece puentes que conectan diversas áreas, convirtiéndose por ello en un elemento que promueve la visión interdisciplinaria de un objeto de estudio. En otras palabras, la estadística se encuentra presente en la educación, la psicología, la antropología, la economía, la sociología, la geografía, etc. Todas estas disciplinas cuando convergen en un tema común, pueden conectarse desde el lenguaje que proporciona el método científico y su principal auxiliar que es la estadística.

Como se dijo, la estadística deriva de las matemáticas y por tanto comparte el lenguaje de los números, pero difieren en la manera en que los números cobran sentido. Buena parte de la estadística puede entenderse desde las operaciones básicas y un álgebra elemental. Lo verdaderamente creativo en la estadística, es cómo tales operaciones se utilizan para responder preguntas concretas y direccionar acciones en la vida real. Lo más importante para aprender en esta disciplina, es el modo en que puede encararse un problema usando una lógica y un razonamiento particular. Una vez que se logra comprender el fundamento de un modelo estadístico, éste forma parte de una manera de pensar la realidad. Permite extender los límites de lo particular y enfocarse en lo general, es factible ver el comportamiento de los sistemas, establecer series de tiempo evolutivas o históricas, y lo más importante, en la investigación científica nos permite evaluar la verosimilitud de hipótesis teóricas e identificar los alcances y límites de las empíricas.

Esta visión de conjunto está plasmada en la etimología de la palabra estadística, la cual deriva de la voz latina “Estado” (*status*). En tal sentido, se sabe que el Imperio Romano llevaba un minucioso compendio de los movimientos poblacionales y las riquezas de las regiones y sus habitantes. Anteriormente, en Grecia, Sócrates hablaba de lo necesario que es para el gobierno el conocimiento de la población y la riqueza de las ciudades. En Egipto, la organización administrativa del estado se basaba en observaciones sistemáticas y periódicas de la población y las producciones agrícolas de las distintas regiones, que se volcaron a tablas que han sobrevivido hasta nuestro tiempo. La estadística aparece en textos bíblicos, tal el caso de Los Números, donde se

cuenta que Moisés ordenó el primer censo de los israelitas dispersos por el desierto. También se sabe del censo ordenado por David, por el relato en el segundo “Libro de Los Reyes”. En China el rey Yao (3000 a.C.), ya contaba con un censo de población y producción agrícola.

La caída del Imperio Romano y su disolución dejó escasos documentos estadísticos y censales, y se produjo una merma notable en la recolección y análisis de información. El rey Guillermo de Inglaterra reconoció la importancia de los censos de población y ordenó que se censara la población en todo su dominio en el año 1000 aproximadamente. La iglesia, en el Concilio de Trento, introduce en forma obligatoria el registro de nacimientos, matrimonios y defunciones. Se reconoce nuevamente la importancia de los registros numéricos de población para las cuestiones estatales, y se debe a German Cönnig las primeras ideas sobre registros y análisis de datos estadísticos, disciplina que llamó Estadística Universitaria. El trabajo de Cönnig es lo que hoy forma el apartado de la estadística descriptiva. En Inglaterra John Graunt comienza a utilizar los registros estadísticos más allá de la simple descripción. Su profesión de demógrafo lo lleva a fundar las bases de la bioestadística y se lo considera el precursor de la Epidemiología. Aunque sus investigaciones estuvieron alejadas del método científico actual, le dio a la estadística un uso que amplió notablemente sus horizontes.

Por aquella época destacados matemáticos se habían interesado en el problema del azar y el cálculo de probabilidades, pero fue Jaques Bernouilli quien introduce los conceptos de certeza y probabilidad en problemas sociales. Los aportes de Abraham de Moivre y Pierre Simon Laplace al cálculo de probabilidades, le permitió a Karl Friedrich Gauss desarrollar su teoría sobre la distribución de los errores y proponer el modelo estadístico más reconocido: la curva de Gauss o distribución normal. Adolfo Quetelet aplica el modelo aportado por Gauss a estudios sociales y antropológicos y se lo considera el padre de la Antropometría. El mismo modelo de distribución normal de Gauss fue utilizado por Francis Galton en sus estudios sobre herencia e inteligencia, y aparece con él la disciplina de la Psicometría. La estadística sigue un curso que se imbrica cada vez más con los modernos métodos de la estimación y la investigación científica, y son hombres interesados en diferentes ramas de la ciencia los que siguen dando impulso a los nuevos métodos de análisis estadísticos. Nombres reconocidos de la moderna estadística son Karl Pearson, Sir Ronald Fisher, George Udny Yule, William Sealy Gosset, Charles Spearman, entre otros.

En este breve compendio histórico puede verse que la estadística tiene siglos de historia y se ha ganado un lugar de preponderancia en muchas disciplinas sociales. En la actualidad y con los modernos métodos de cálculo apoyados en el uso de computadoras, pueden obtenerse indicadores estadísticos casi al instante. Pero estos

cálculos por si mismos no tienen demasiado sentido si no hay personas preparadas para darles un uso racional. Un curso de estadística, por breve que sea, debe introducir no solo las nociones básicas; también debe fortalecer la capacidad analítica del estudiante. Actualmente los informes estadísticos se encuentran disponibles en publicaciones con diversos soportes. La alfabetización estadística es ya una herramienta del pensamiento profesional que deben dominar los estudiantes de humanidades, especialmente aquellos que se interesan por la educación, puesto que los sistemas internacionales de indicadores educativos son estadísticos.

Como ya se mencionó, la Estadística es una rama de las matemáticas que provee técnicas para el trabajo con datos, dichas técnicas surgen en su mayor parte de la necesidad de fundamentar la toma de decisiones en condiciones donde prevalece la incertidumbre. La noción de incertidumbre estadística se relaciona con la idea vulgar de indecisión de la vida cotidiana, pero se diferencia fundamentalmente en que bajo los modelos estadísticos, la incertidumbre es un aspecto mensurable y por tanto conocido, de la realidad. Por ejemplo, al leer textos informativos, académicos o de divulgación, aparecen las nociones estadísticas sobre las cuales se basan los juicios valorativos. Sirvámonos de un ejemplo: la versión electrónica del diario El Cronista, publica el 25 de Abril de 2016¹ un informe del Ministerio de Seguridad sobre el mapa de delitos en la República Argentina, destacándose que hubo una caída del delito entre los años 2014 a 2015, aunque las tasas se mantienen por encima de las de 2008. Refiriéndose específicamente a delitos contra las personas y la propiedad, se subraya que la tasa nacional de este tipo de delitos cayó un 3% entre los años mencionados aunque se mantiene un 10,2% por encima del 2008. Vemos en este escueto ejemplo la manera en que se intenta responder a la pregunta sobre la inseguridad utilizando valores porcentuales para mensurar el impacto sobre la sociedad. Las interpretaciones de estas cifras de seguro van a cambiar de un analista a otro, pero esa discusión excede este trabajo. Lo que destacamos es que la estadística construye sus datos, y en este sentido el interés del analista o investigador no puede prescindir del pensamiento teórico. En otras palabras, la estadística por sí misma no es suficiente para comprender cabalmente los problemas de los que se ocupa. Son los investigadores los que se valen de ella para enfrentar los problemas y tomar medidas. Al construir sus datos, la estadística define un lenguaje propio de la disciplina. Así, términos como variables, unidad de análisis, medición, variabilidad, estimación y otros, se conjugan en lo que constituye la semántica propia del análisis estadístico. Es necesario que los estudiantes comprendan paulatinamente que las afirmaciones de tipo estadístico, tales como chance o probabilidad, se ajustan a modelos matemáticos que los investigadores utilizan para describir la realidad, y que se apartan del sentido común dado a esos términos.

Los datos que se recogen siempre refieren a unidades de análisis y los

estadísticos calculados son formas de caracterizar a ellas. De este modo, las tasas de repitencia de un determinado nivel del sistema educativo, refiere a unidades de análisis identificadas como alumnos. Pero el foco de este análisis está puesto sobre una propiedad que trasciende al alumno y sirve para caracterizar un colectivo o grupo agregado de unidades individuales. Dicho en otras palabras, si un alumno no es promovido al año siguiente superior, se cuenta entre aquellos que repiten el grado. Del conteo de estos alumnos surge una cifra que es la cantidad de repitentes, que al ser ponderada por la cantidad total de alumnos se transforma en una proporción que se interpreta como tasa de repitencia. Este indicador, nos habla de la situación de un colegio, una jurisdicción, una provincia o el país. Las técnicas estadísticas en su posibilidad de agregar y desagregar datos nos muestran el comportamiento y dinámica de variables propias de una población (estudiantes para este ejemplo). El potencial del análisis estadístico se aprecia más cuando es posible combinar sus métodos con otras formas de análisis. En educación es común que se hable de un universo “micro” si el interés recae en la escuela, el grado o el alumno, luego el universo “macro” caracteriza los colectivos agregados como los alumnos o las escuelas. La frontera entre estos términos es arbitraria y subjetiva; de modo que los investigadores siempre tienen que hacer explícitas sus técnicas de análisis y es en ese instante donde pueden darse intersecciones que operen rupturas con los velos que impone un saber cotidiano y no cuestionado, lo dado y aceptado, en fin, con los prejuicios.

Este tutorial ha sido redactado de manera que pueda ser usado como material introductorio a los principales temas de la estadística. Se verá aquí que el contenido matemático ha sido reducido al mínimo, aunque no se prescinde de algunas fórmulas conceptuales y ciertas rutinas de cálculo. El uso de estas fórmulas tiene por finalidad que los estudiantes comprendan los conceptos que subyacen al estadístico que se calcula, mediante la simbología matemática.

Capítulo 1

Variables y Unidad de Análisis

Las variables son la piedra angular de un análisis estadístico. Constituyen básicamente lo que el investigador resalta como importante en un estudio, y determina el modelo de análisis que podrá aplicar para encontrar patrones dentro de los datos, someter a prueba hipótesis, estimar parámetros, etc. El ajuste de estos modelos depende además de cómo se haya realizado la medición, de ello resulta que existen diferentes sistemas de medición aplicables a distintos tipos de variables. Finalmente, las variables pueden seleccionarse por los atributos directamente observables en la unidad de análisis, o pueden definirse teóricamente. De esto último deriva la idea de rasgo latente, que identifica aquellas variables que guardan un isomorfismo entre el atributo y la medición.

En este capítulo abordaremos el tema de variables, unidad de análisis y sistemas de medición. Estos tres tópicos conforman un conjunto que se definen mutuamente y que son la base para el análisis y resumen de datos. En otras palabras, dependiendo de cuáles sean las variables de interés, podremos escoger el sistema de medición más conveniente a los objetivos de un estudio. Determinar qué vamos a medir y cómo lo mediremos, implica definir la unidad de análisis correspondiente. Como veremos aquí, Variable – Unidad de Análisis – Sistema de Medición, son tres elementos fundamentales para cualquier operación con datos. En estadística es importante saber cuál es el dato que estamos buscando para cumplir con el propósito de un estudio. El dato, que es el pilar fundamental donde se asienta el análisis de cualquier problemática, depende de lo que hemos definido previamente como variable. La variable, en términos sencillos, nos indica qué es lo que estamos midiendo y la unidad de análisis nos indica donde medirlo.

Hay que tener en cuenta que la variable importa en la medida en que luego nos indica algo acerca de la unidad de análisis, por lo que ambas no pueden existir de manera independiente. Las técnicas y modelos estadísticos son diferentes según el tipo de datos con que se cuente, es decir, la variable que se analiza. Existen diferentes clasificaciones para los datos, pero aquí intentaremos simplificarlas. En principio los datos pueden clasificarse en dos grandes categorías que son:

- a) datos cuantitativos, que consisten en números que resultan de conteos o mediciones,
- b) datos cualitativos que consisten en atributos o categorías a las que pueden

asignársele valor numérico, pero no puede operarse matemáticamente con ellas.

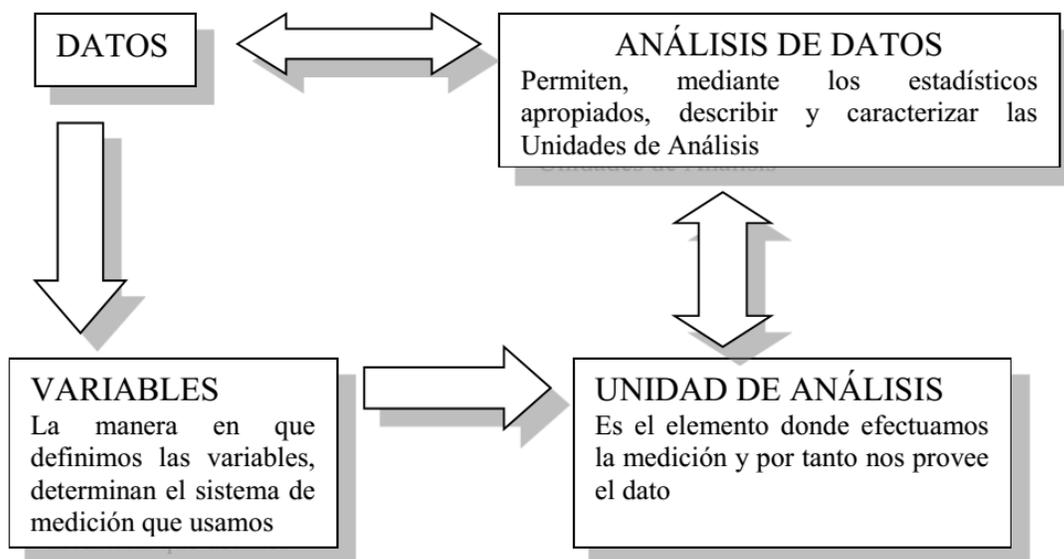
En este punto conviene discriminar los términos cuantitativo y cualitativo cuando son aplicados a los datos, de las veces que se usan para referirse a la metodología de la investigación. Veamos la diferencia utilizando un ejemplo ficticio. Supongamos que estamos relevando información de la cantidad de profesores con título de Doctor en diferentes disciplinas; en este estudio, también recogemos datos de la edad a la que se doctoró la persona. La edad es un dato cuantitativo, y como tal es tabulado numéricamente. La disciplina en la que recibió el título de Doctor, es un dato cualitativo y será tabulado mediante una etiqueta, por ejemplo, Doctor en Filosofía, Doctor en Letras, Doctor en Ciencias de la Educación, etc. Promediando las edades y ordenándolas por disciplinas, podríamos obtener la siguiente tabla:

Disciplina en la que se doctoró y edad promedio de obtención del título máximo

<i>Disciplina</i>	<i>Edad Promedio (en años)</i>
Filosofía	34
Letras	36
Cs. Educación	37
Historia	41
Cs. Sociales	42
Psicología	44

En la tabla se observa que los profesores que obtienen su grado de doctor a menor edad, pertenecen en su mayoría a la disciplina Filosofía, mientras que los que se gradúan más tardíamente con el título máximo, tienden a estar agrupados en la disciplina Psicología. La tabla también separa el dato cualitativo en la primera columna, y el dato cuantitativo en la segunda. Para comprender mejor la interrelación entre datos, variables y unidad de análisis, presentamos el esquema que sigue a continuación:

Datos, Variables y Unidades de Análisis



Veamos ahora con unos sencillos ejemplos la forma en que variable y unidad de análisis quedan identificadas en el contexto de un problema.

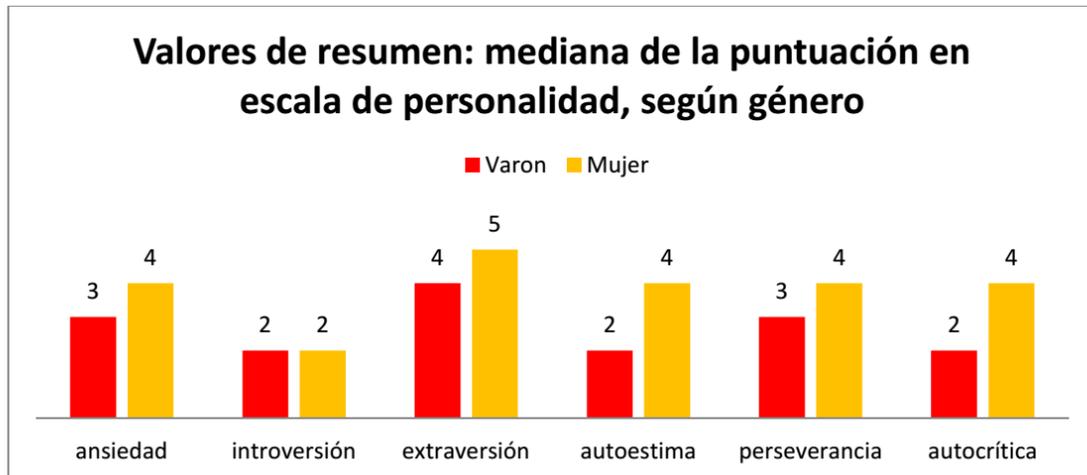
Un investigador se encuentra realizando un estudio sobre personalidad en estudiantes universitarios de la escuela de enfermería. Para ello, selecciona una muestra accidental de 100 estudiantes (50 mujeres y 50 varones) de esa carrera, a quienes aplica un cuestionario de personalidad, cuyos resultados se expresan de manera cuantitativa con puntajes de 1 a 5 (1=poco; 5=mucho) en los siguientes factores: a) ansiedad, b) introversión, c) extraversión, d) autoestima, e) perseverancia, f) autocrítica. El estudio pretende determinar si existe un perfil de personalidad diferente según el género de los estudiantes.

Primero identificamos la variable que está en estudio, en este caso se trata de un perfil de personalidad definida por los factores descritos desde a) hasta f). Esta variable se mide con un instrumento que es un cuestionario; el cual se aplica a una muestra de estudiantes de la carrera de enfermería, diferenciados según el género. Por lo tanto tenemos que:

Unidad de Análisis: estudiantes de la carrera de enfermería.

Variables: a) Género del estudiante, b) Factores del cuestionario de personalidad.

En este ejemplo ficticio, podríamos obtener el perfil de personalidad diferenciado por género, tomando como estadístico de resumen la mediana de las puntuaciones en el cuestionario de personalidad, lo cual podría resumirse en el siguiente gráfico.



Como se aprecia en el gráfico, las mujeres tienden a mostrarse más ansiosas que los varones. En cuanto a la introversión, no hay diferencias por género y la tendencia muestra que no se perciben introvertidos. En extraversión se ve que ambos sexos se perciben muy extrovertidos, aventajando las mujeres a los varones. En autoestima y autocrítica, las mujeres aventajan a los varones, percibiéndose con autoestima alta y mucha autocrítica. En perseverancia, ambos sexos se ven perseverantes, con una ventaja de las mujeres sobre los varones.

En este punto no llevaremos más lejos las interpretaciones dado que no se trata de todavía no hemos visto medidas de tendencia central y muestreo, además no se trata de una investigación real. El ejemplo pretende hacer notar que una vez que hemos definido la variable, su escala de medición y hemos obtenido los datos, éstos se resumen con un estadístico, y la información se presenta en un gráfico desde el cuál hacemos la lectura de los mismos. Es importante notar desde ya, que la conclusión aporta información sobre las unidades de análisis, y esa información está resumida mediante estadísticos.

Sistemas de Medición

Siempre que definimos una variable de interés para un estudio, es porque vamos a efectuar una medición sobre ella. La definición de medición que vamos a dar aquí es de tipo general, de manera que nos permita entender que existen diferentes sistemas de medición cada uno con características particulares. Entonces, efectuamos una

medición cada vez que asignamos un valor a una variable siguiendo una regla. Debe notarse aquí que la palabra valor, no alude necesariamente a un valor numérico y que es el propio investigador quien puede determinar cuál es la regla de asignación. Una vez establecida la regla, no puede ser cambiada en el curso de la recogida de datos, por lo cual debe asegurarse que la regla pueda contener a todas las unidades de la muestra o población. Es decir, el modo en que asignamos valor a la variable en estudio debe abarcar a todas las unidades de análisis.

Sistema Nominal: el sistema de medición nominal es el más simple y extendido de los sistemas de medición, consiste en crear categorías para clasificar a las unidades de análisis. Este sistema puede admitir valores numéricos, pero solo a título de etiquetas. El ejemplo más sencillo para captar el nivel de medición nominal, son las categorías SI – NO – NS/NC. Por ejemplo, se podría hacer un sondeo de opinión electoral con la siguiente pregunta: *¿votará Ud. por el candidato Pérez en las próximas elecciones?* Podríamos usar las categorías de respuestas para clasificar a los encuestados, o suplantarnos por números de la siguiente manera: 1= SI; 2= NO; 3= NS/NC. Nótese que los números solo designan el tipo de respuesta, y no es lícito suponer que la respuesta SI es menor que la respuesta NO, por tener ésta última el valor 2. Suponiendo que existan en la contienda electoral tres candidatos con la mejor posición en comparación con el resto, se podría crear una escala nominal con las siguientes categorías: Pérez, López, Gutiérrez, Otro. En este caso la pregunta en el sondeo de opinión sería *¿Por cuál candidato votará en las próximas elecciones?* Si la elección de la persona encuestada no se inclina por los primeros tres candidatos, su respuesta se clasificará en la categoría Otros. Esto último es una regla de asignación que deberá ser respetada cada vez que se clasifique la respuesta del encuestado.

Muchas características de las personas tienen implícito un sistema de clasificación nominal, por ejemplo: a) lugar de procedencia, b) profesión, c) género, d) situación laboral, etc. Sugerimos al estudiante consultar la encuesta permanente de hogares (EPH), donde se presentan muchas variables de este tipo y su regla de asignación (la EPH es realizada con periodicidad por el INDEC y puede consultarse en el siguiente enlace <http://www.indec.gov.ar/>). Un sistema de clasificación nominal exige casi siempre que las categorías sean exhaustivas y mutuamente excluyentes. Es decir, cada vez que usamos este sistema se espera que la persona quede registrada solo en un valor de la escala; por ejemplo, no se espera que un individuo en la variable situación laboral, sea clasificado al mismo tiempo como empleado y desempleado. De repetirse esto, generaría un doble conteo de unidades de análisis y la base de datos contendría muchas más respuestas que individuos. Sin embargo, se puede encontrar situaciones en las que es necesario clasificar las respuestas como categorías. En este caso, se dice que el sistema clasificatorio es múltiple, y lo que se toma como unidad de

análisis son las respuestas dadas y no a los individuos. Un ejemplo de ello es cuando se pregunta en sondeos de opinión cuál o cuáles son las preferencias por objetos, lugares o consumos.

Siguiendo este razonamiento, supongamos que preguntamos a una muestra aleatoria de 200 estudiantes de la Facultad de Filosofía y Humanidades, cuáles son las actividades recreativas preferidas los fines de semana. Nótese que en principio no hay categorías previas de clasificación, sino que estas quedarán establecidas a partir de las respuestas recogidas. De este modo, un individuo puede optar por el deporte y la vida al aire libre, mientras que otro puede optar por el deporte, el cine y los recitales, otro incluso podría optar por un paseo de compras, etc. Si bien este sistema de clasificación es nominal, se diferencia de una escala nominal en que: a) se clasifican las respuestas y no los individuos, b) pueden no tener un sistema clasificatorio establecido de antemano, c) los individuos pueden dar una o más respuestas, e) es exhaustivo, pero las respuestas no son mutuamente excluyentes. Para sintetizar diremos que si tomamos a los individuos como unidad de análisis, las categorías de las variables son mutuamente excluyentes; en cambio, cuando la unidad de análisis son las respuestas de los individuos no se satisface necesariamente esta condición.

Sistema Ordinal: este sistema es idéntico al anterior salvo en que los valores preservan una relación de orden, aunque su naturaleza no es métrica. Dicho en otros términos, en un nivel de medición ordinal, los datos se agrupan de acuerdo a un orden expresado con los operadores “más que”, “mayor que” (o su contrario implícito “menos que”, “menor que”). Los sistemas ordinales son muy utilizados en los casos en que deben establecerse juicios de orden en la variable, pero no es posible establecer iguales distancias entre las valoraciones. El nivel de medición ordinal también puede usarse para asignar rangos a las categorías, práctica muy común en las ciencias sociales. Los rangos ordenados pueden extender el número de categorías y aproximarse a una escala métrica. Veamos un ejemplo: se intenta conocer la opinión de los profesores de una escuela sobre la calidad de la gestión del cuerpo directivo de la misma. Para ello se realiza una encuesta anónima pidiendo la valoración de dicha gestión utilizando el rango 1 a 10, donde 1 indica una mala gestión y 10 una excelente gestión. La afirmación que deberán puntuar los profesores sería la siguiente: *“Valore de 1 a 10 en qué medida considera provechosa para la escuela la actual gestión del cuerpo directivo”*. Cabe destacar que muchas encuestas usan este tipo de afirmaciones como disparadores para valorar las respuestas de las personas, dichas afirmaciones se denominan “reactivos”, que no es otra cosa que la oración a la que deberá otorgar un puntaje el encuestado. Técnicamente, cada pregunta, afirmación o ítem de un test o cuestionario, recibe el nombre de reactivo en tanto es el estímulo al cual debe responder el individuo. Continuemos con el ejemplo y supongamos que respondieron

100 docentes. Es necesario ahora reordenar los datos para describir el resultado obtenido. Esto puede hacerse mediante una tabla de frecuencia que indique en que porcentaje los docentes han repartido sus respuestas. El resultado quedaría expresado de la siguiente manera.

Distribución porcentual de la valoración de la gestión escolar

1 - 2	7%
3 - 4	12%
5 - 6	9%
7 - 8	44%
9 - 10	28%

Como se aprecia, los valores de la tabla han sido agrupado de a dos para reducir su tamaño, además se ve que la gestión de la escuela tiene una alta valoración en la mayoría de los docentes que respondieron. No avanzaremos más en la interpretación, puesto que tablas de frecuencia es tema que abordaremos en el capítulo siguiente.

La asignación de rangos puede aplicarse a casos en que el atributo medido pueda oscilar en dos extremos opuestos. Imaginemos la siguiente situación: se propone una escala de 1 a 5 para valorar el estado de felicidad de una persona. El reactivo usado podría ser el siguiente: *“Valore de 1 a 5 en qué medida Usted se siente feliz en este momento”*. La única conclusión que podríamos sacar de aquellas personas que se categorizan en el valor 1, es que no están felices, pero no podríamos afirmar que están tristes, pues esa información no es captada ni por la escala, ni por el reactivo. Para poder llegar a esa información, será necesario modificar ambos. Primero agregaríamos una categoría neutra, el cero, y luego valores positivos y negativos. La ponderación quedará expresada en una escala como la que se muestra a continuación:

-5 -4 -3 -2 -1 0 +1 +2 +3 +4 +5

Ahora deberíamos modificar el reactivo de la siguiente manera: *“Valore su estado de ánimo en este momento, siendo 0=neutro; -5 muy triste; +5 muy feliz”*. Los rangos ordenados nos proporcionan ahora mucha más información, y sabremos que aquellos individuos que se categorizan en el extremo positivo de la escala son quienes experimentan un estado de ánimo feliz, al contrario, los que se categorizan en el extremo negativo manifiestan un estado de ánimo triste. El cero es el punto de referencia de la escala de rangos ordenados. Algo que no debemos olvidar en este tipo de escalas, es que las distancias entre los valores no tienen intervalos iguales, por lo tanto una persona que puntúa 4 en la escala anterior, no está doblemente feliz en comparación con otra persona que puntúa 2. Simplemente diremos que el primer individuo está “más feliz” que el segundo.

En ciencias sociales un aporte fundamental a este tipo de escalas de medición

fue realizado por Rensis Likert, quien fue educador y psicólogo organizacional. Como parte de su trabajo en este último campo desarrolló la metodología de medición que se conoce como escala de Likert. Este tipo de escala posee los atributos que hemos mencionado anteriormente, y resultan fundamentales para valorar objetivamente cuestiones que tienen un carácter subjetivo. Rensis Likert desarrollo escalas basándose en la asignación de valores numéricos para medir las motivaciones, los estilos de liderazgo, las preferencias, los acuerdos, el agrado o desagrado con ciertas afirmaciones, las inclinaciones por un método, modelo o candidato, las intenciones, etc. Actualmente se reconoce su aporte en que cada vez que se crea una escala ordinal con asignación numérica a la valoración, se la llama escala tipo Likert.

Hasta ahora hemos descripto escalas ordinales numéricamente valoradas, pero esta propiedad no es una condición necesaria. Esto quiere decir que podemos utilizar etiquetas de valor para confeccionarlas. Las etiquetas indicarán en cada caso la distancia entre las categorías. Veamos un ejemplo: se consulta a una muestra aleatoria de 200 docentes de diferentes Facultades de la UNC, sobre la implementación del examen de ingreso obligatorio. La valoración se realiza para la siguiente afirmación: *“en qué medida acuerda usted con el examen de ingreso a la Universidad Nacional”*. Los resultados de esta hipotética encuesta se resumen en la siguiente tabla:

Nivel de acuerdo de docentes universitarios con el examen de ingreso

Acuerdo completamente	11%
Acuerdo moderadamente	4%
Desacuerdo moderadamente	38%
Desacuerdo completamente	47%

Las cuatro etiquetas de valor pueden ser reemplazadas por números, pero para este ejemplo vemos que conviene más conservarlas, dado que describen mejor el resultado obtenido. Por otro lado, aunque no se utilicen números se observa que las categorías están ordenadas en la medida en que establecen un continuo entre aquellos que están completamente de acuerdo con el examen de ingreso, y quienes desacuerdan completamente con el mismo. La interpretación de este resultado indicaría que la mayoría de los docentes se manifiesta en desacuerdo con el examen de ingreso a la Universidad Nacional.

Existe una regla empírica a la hora de crear una escala ordinal para una investigación, y es que si la escala puede reducirse a unas pocas categorías (no más de cinco), conviene utilizar etiquetas de valor. Si la escala debe captar más información y

por tanto se necesitan más categorías, conviene valorarla numéricamente. Esta regla deriva del hecho de que es más fácil utilizar descriptores verbales cuando las categorías son pocas, situación que se complica notablemente si se usan muchas categorías.

Sistema Métrico: en este sistema además de establecerse una relación de orden, se verifica que las diferencias entre las magnitudes se corresponden puntualmente con las diferencias de los objetos de medición. Esto permite establecer diferencias numéricamente ponderadas entre objetos; un ejemplo al que estamos acostumbrados y que ilustra lo antes dicho, son las diferencias en litros en el contenido de las bebidas. Así, dos gaseosas de 1,5 litros tienen el mismo contenido que tres gaseosas de 1 litro. Todos los patrones de medida son sistemas métricos, y solo en este caso es posible operar matemáticamente con las magnitudes obtenidas. Dentro de los sistemas métricos pueden distinguirse dos tipos de escalas:

Escalas Intervalares: en un nivel de medición intervalar los datos se ordenan en intervalos iguales; aunque se puede operar matemáticamente con los valores de la escala, no es posible establecer razones dado que no se cuenta con un cero absoluto. El mejor ejemplo para ilustrar este nivel de medición son las escalas de temperatura. En nuestro país es de uso corriente la escala Celsius, cuyo cero está dado por el punto de congelación del agua y el 100 por su punto de ebullición (de allí que se denomine escala centígrada). Se comprende entonces que decir cero grado no implica ausencia de temperatura; es simplemente el referente empírico en donde se sitúa el valor cero. El sistema centígrado no es el único sistema para medir temperatura, y de hecho en el sistema internacional la temperatura se mide en grados Kelvin. El sistema centígrado y el sistema Kelvin emplean intervalos iguales cada uno, pero el tamaño del intervalo no es el mismo en las dos escalas y el punto cero tampoco. Por lo tanto, para establecer una relación entre ambas escalas, es necesario emplear una ecuación que represente la transformación entre escalas. En otras palabras, estos dos sistemas son intercambiables mediante la siguiente ecuación:

$$C^{\circ} = K^{\circ} - 273.15$$

$$K^{\circ} = C^{\circ} + 273.15$$

Como ya mencionamos, lo que diferencia una escala intervalar de una proporcional es que el punto cero es arbitrario. Dado este hecho es incorrecto hablar de razón en estas escalas, lo cual puede mostrarse fácilmente con el siguiente ejemplo. Imaginemos que tenemos una escala graduada del 0 al 9 para medir algo, y que por comodidad desplazamos el cero una unidad, tendríamos una situación como la siguiente:

0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9

En la escala superior entenderíamos que 4 es el doble de 2 y que 8 es el doble de 4. Al ser escalas de intervalos iguales, podríamos proyectar la relación de la escala superior en la inferior, y se produciría una inconsistencia ya que no se cumple que 5 sea el doble de tres ni nueve el doble de cinco. Es decir, al ser el origen arbitrario, no es posible equiparar sin más ambas escalas, aun cuando posean intervalos iguales.

¿Es un problema no contar con un punto de origen fijo? En estadística esto no es realmente un problema, pues las escalas de intervalo permiten operar matemáticamente y por tanto es posible calcular todos los estadísticos. Además, muchas escalas de medición de uso común, tienen la propiedad de intervalos iguales y cero arbitrario, por ejemplo la altura desde el nivel del mar, la distancia de una localidad respecto de un punto de referencia, el índice de masa corporal, etc. Por ello, los paquetes informáticos actuales no hacen distinción entre las escalas de intervalo y de razón, y simplemente se refieren a ellas como escalas métricas.

Escalas Proporcionales o de Razón: este nivel de medición es similar al nivel anterior, pero con la propiedad de que el cero tiene valor absoluto. Es decir, allí donde se registre cero, será igual a nada en la variable. Durante mucho tiempo se ha trabajado para que todas las unidades del sistema internacional estén expresadas en escala de razón. Aunque no se ha logrado completamente, se avanzó mucho en este sentido; así la longitud se mide en metros, la masa en kilogramos, el tiempo en segundos, etc. De este modo, un elemento que haya recorrido 0 metro, indica que no se ha movido de lugar, un elemento con 0 kg de masa, no posee masa, y el instante 0 segundo, es el momento inicial.

En dos escalas de razón existe la misma relación entre los intervalos y además existe la misma relación entre dos puntos de la escala. Por esta propiedad, las transformaciones entre escalas de razón son monotónicas y quedan expresadas por una función, que en notación matemática es $y = f(x)$. Esto quiere decir que cualquier valor de y puede calcularse a partir de los valores de x , mediante la función f . Si x e y son dos escalas de razón, tendremos correspondencias entre ambas para cualquier valor dado.

Como ya mencionamos, en estadística no es tan importante la distinción entre escalas de intervalo o de razón, ya que ambas serán tratadas como escalas métricas. Lo importante es recordar que solo en este tipo de escalas pueden emplearse operaciones matemáticas para los estadísticos que las requieran, como por ejemplo el promedio o la varianza. Asimismo, en ciencias sociales las escalas métricas se emplean

cuando: a) es posible medir con un patrón, el caso de la altura o el peso de una persona; b) se pueden realizar conteos, por ejemplo, cantidad de materias aprobadas, cantidad de alumnos inscriptos al comienzo del año lectivo, y c) cuando se emplean pruebas estandarizadas, por ejemplo el puntaje de la evaluación del Programa Internacional para la Evaluación de Estudiantes (PISA por sus siglas en inglés: *Programme for International Student Assessment*). Otra cuestión a destacar es que las escalas métricas pueden reducirse convenientemente a escalas ordinales, pero una escala ordinal no podrá ser tratada como métrica. Para ilustrar este punto tomemos un ejemplo sencillo. Supongamos que medimos la altura de cien personas. Cada medición quedará expresada por un valor en centímetros y así habremos obtenido cien mediciones que corresponden a una escala métrica. Podría darse el caso de que las cien mediciones efectuadas fueran diferentes, aunque tales diferencias fueran mínimas; por ejemplo: el individuo A tiene una altura de 1,71 cm, el individuo B tiene una altura de 1,65 y el individuo C una altura de 1,66. Bajo estas circunstancias, podríamos reducir la escala métrica original a una escala ordinal con las siguientes categorías; a) Persona de estatura baja: 1,60 cm o menos; b) Persona de estatura media: 1,61 cm a 1,70 cm; y c) Persona de estatura alta: 1,71 cm o más. La escala métrica original quedó expresada en una escala ordinal con tres categorías: Persona de estatura Alta – Media – Baja. Aunque ahora hay menos categoría y la escala es más simple de manipular, se ha perdido información de la escala métrica original, pues si decimos que el individuo D es de estatura Alta, mientras que el individuo E es de estatura Media, solo sabremos que D mide 1,71cm o más y que E tiene una estatura comprendida entre 1,61 y 1,70 cm. Si las cien personas hubieran sido clasificadas en una escala ordinal como la que creamos, sería imposible reconstruir la escala métrica original, y esto es lo que se quiere expresar cuando afirmamos que una escala ordinal no podrá ser tratada como métrica.

Otro uso del término escala

La clasificación de las escalas de medición que hemos ofrecido en los párrafos anteriores, corresponde a la propuesta de S. Stevens, quien en 1946 publicó un artículo de amplia difusión en la revista *Science* (Nº103, pp: 677-680), titulado *On the theory of Scales Measurement* (Sobre la teoría de las escalas de medición). Desde entonces se ha mantenido por conveniencia la diferenciación presentada. Sin embargo, en la moderna investigación social, han aparecido escalas de medición que no fueron contempladas en ese artículo. Mencionaremos solo tres que son las más conocidas: a) Escala porcentual: el porcentaje consiste en tomar cualquier valor y reducirlo a una base cien, de allí que sea posible establecer cambios porcentuales, incrementos o decrementos en los mismos. De esto se derivó que los porcentajes puedan ser tratados como escalas de medición de tipo métricas y calcularse diferentes estadísticos a partir

de los mismos, de hecho, el promedio para una escala porcentual se denomina media geométrica; b) Escala autograduada: esta fue una invención muy útil en la medición del dolor, cuestión que durante mucho tiempo no contó con un procedimiento metódico. La escala fue creada combinando las propiedades de la escala de intervalo y ordinal. El procedimiento consistió en darle al paciente una regla graduada en centímetros del 0 al 10; se le informaba que cero es ausencia del dolor y diez es un dolor imposible de soportar sin analgésicos. Los pacientes se familiarizaban con el tamaño de la regla, pero luego debían sostenerla de manera tal que no pudieran ver los números. Así, debían indicar cuánto dolor experimentaban, señalando con el dedo en el dorso de la regla. Sorprendentemente se encontró en los reportes que el umbral del dolor soportable sin analgesia estaba entre el 6 y el 7. La familiaridad de las personas con una regla de 10 cm., hizo que su ponderación del dolor se volviera regular. A partir de este hecho, muchas escalas han aprovechado la propiedad de mediciones conocidas por las personas para ponderar aspectos muy subjetivos como por ejemplo las preferencias estéticas; c) Pruebas estandarizadas: el auge de las pruebas estandarizadas se dio entre el fin del siglo XIX y principios del siglo XX. La escala estandarizada más conocida es la desarrollada en 1904 por Alfred Binet para medir la inteligencia. La teoría de la medición psicométrica, encontró por aquella época muchos seguidores y su método pronto se extendió a otros campos de las ciencias humanas. Los aportes fundamentales de Chales Spearman, Edward Thorndike, Lewis Therman y Robert Woodworth, sirvieron para desarrollar numerosas escalas para medir no solo aspectos psicológicos, sino también factores sociales y antropológicos. Existen dos cuestiones fundamentales que definen a una escala estandarizada; la primera es que se cuenta con valores de referencia para una población o grupo etario. La segunda es su validez y confiabilidad. El primer aspecto se refiere a que, al utilizar una prueba estandarizada en un conjunto de individuos, resulta factible realizar una comparación con un patrón de referencia, cuyo nombre técnico es baremo. La segunda cuestión se refiere a que la prueba ha demostrado que efectivamente mide lo que se propone medir, y que lo hace de manera consistente. Más adelante discutiremos estos puntos en el apartado variables de constructo. Por el momento basta con recordar que las pruebas estandarizadas contienen en su mayoría escalas de tipo métricas.

Capítulo 2

Técnicas de muestreo

La estadística cuenta con diversas herramientas para estimar en una población los parámetros que deseamos conocer. La estimación dependerá en buena medida de una apropiada definición operacional de lo que entendemos como población. Definir una población con precisión implica un trabajo detallado y minucioso de sus límites, de modo que sea factible obtener una muestra representativa de la misma. Si estamos interesados en estimar un parámetro, tendremos que optar por un sistema de muestreo que garantice la confianza en la estimación. Es por ello que apelamos a las técnicas de muestreo probabilístico. Junto al muestreo probabilístico, existen otras técnicas menos rigurosas para obtener muestras de una población. Aunque no se emplean en la estimación de parámetros, los tipos de muestreo no probabilístico se aplican en casos puntuales aportando valiosa información del comportamiento de una variable de interés. En este apartado veremos la manera en que se define teórica y operacionalmente una población, y cómo se aplican las distintas técnicas de muestreo, tanto probabilístico como no probabilístico.

Muestra y Población

Supongamos que ponemos a hervir una olla con arroz; pasado cierto tiempo deseamos saber en qué punto de cocción se encuentra. Entonces, revolvemos el contenido y tomamos una cucharada. De acuerdo al estado de cocción del arroz que tenemos en la cuchara, inferimos que el arroz que se encuentra en la olla debe tener más o menos un estado de cocción similar. Si bien en estadística raras veces nos vamos a ocupar de este tipo de problemas, este simple ejemplo nos ilustra acerca de tres aspectos importantes: a) Población, b) Muestra, c) Inferencia. Si consideramos la olla de arroz como la población de interés, la cucharada extraída correspondería a una muestra de la misma. Al determinar que el arroz de la cuchara se encuentra en un punto de cocción óptimo, inferimos que el resto de la olla debería estar en el mismo punto.

Población

En el lenguaje cotidiano la palabra población nos recuerda a individuos viviendo en una misma región, pero en estadística todo conjunto completo de elementos es una población. Así, todos aquellos individuos que habitan en suelo argentino, conforman la población de Argentina. Podemos extender los límites de este conjunto y decir que la Argentina es un país de Sudamérica y por tanto, es un subconjunto de una población

mayor, la de todos los habitantes de América del Sur. Sin embargo, las poblaciones de las que se ocupa la estadística van mucho más allá de las personas. Podemos definir como población cualquier entidad que pueda formar un conjunto. Por ejemplo, la cantidad de escuelas municipales de la ciudad de Córdoba también corresponde a un conjunto, y puede considerarse una población.

En estadística se dice que una población es finita cuando pueden enumerarse o listarse todos los elementos que forman el conjunto. Por ejemplo, si se define como población a todas las personas que se encuentran en condiciones de votar en la República Argentina, se tiene que la población está listada en lo que se conoce como padrón electoral. Si no podemos establecer el número total de los elementos que contiene la población, decimos que es infinita. Dado que el término finito e infinito se presta a confusión, se está reemplazando paulatinamente por el de límite de la población. Entonces, habrá poblaciones que tendrán un límite acotado, mientras que otras tendrán un límite no acotado. En estadística el límite de la población se toma en un sentido práctico y siempre acorde a los objetivos de la investigación. En este sentido, la primera cuestión que debe resolver el investigador es si tendrá acceso a la población que ha definido teóricamente. Se desprende de ello que la definición de la población determinará el tipo de muestreo posible. Los sistemas de información estadísticos son las fuentes más consultadas cuando se quiere definir el límite de poblaciones a estudiar, especialmente por los datos que allí se recogen sobre población e instituciones. Por ejemplo, si estamos interesados en estudiar algún aspecto relevante de las escuelas municipales de Córdoba Capital, nuestra principal fuente de datos será el Ministerio de Educación. En esa institución están listadas todas las escuelas que conformarían la población de interés.

Existen poblaciones que permanecen no acotadas a pesar de los esfuerzos que se hacen por acotarlas, son las llamadas poblaciones de fracción desconocida. Este es un término técnico, pero con un ejemplo puede entenderse de que se trata. Imaginemos que estamos relevando en Córdoba Capital, la cantidad de personas que padecen de una enfermedad cuya incidencia en la población es muy baja. En primer lugar consultamos los datos del ministerio de salud de la provincia para saber cuántos casos existen en la actualidad, y ese sería el límite de la población. Pero debemos suponer que hay más casos dado que muchos de ellos todavía no han sido diagnosticados. En criminología, se habla de una cifra negra de la criminalidad en tanto se sabe que muchos hechos delictivos no se denuncian, por ende, las estadísticas solo alcanzan a cubrir aquellos casos que efectivamente han sido denunciados. Aquí las estadísticas continuas son de una invaluable ayuda, dado que es posible estimar y hacer proyecciones a futuro a partir de datos conocidos. Tomemos como ejemplo algunas situaciones donde las estadísticas son inexactas pero de inestimable valor: cantidad de casos de aborto, uso de drogas ilegales, abusos infantiles, violencia

doméstica, etc.

Muestra e Inferencia

Cada vez que se seleccione una muestra o una población, está implícito que se desea medir o conocer alguna característica o propiedad importante; ya hemos visto que a estas propiedades las denominamos variables. La pregunta recae la mayoría de las veces en conocer cómo se comportan esas variables en la población. La inferencia es un procedimiento mediante el cual se extraen conclusiones de una población, a partir de una muestra representativa de la misma. Los modelos estadísticos utilizados para el análisis de datos, nos enseñan que una muestra representativa de una población es suficiente para conocer las principales características de ella. Entonces, al utilizar una muestra para determinar una propiedad de la población, estamos realizando una inferencia. Hay una rama de la estadística que recibe el nombre de estadística inferencial, justamente porque desarrolla y aplica métodos para que el proceso de extrapolación de la información que se obtiene de una muestra, se ajuste lo más fielmente posible a la población.

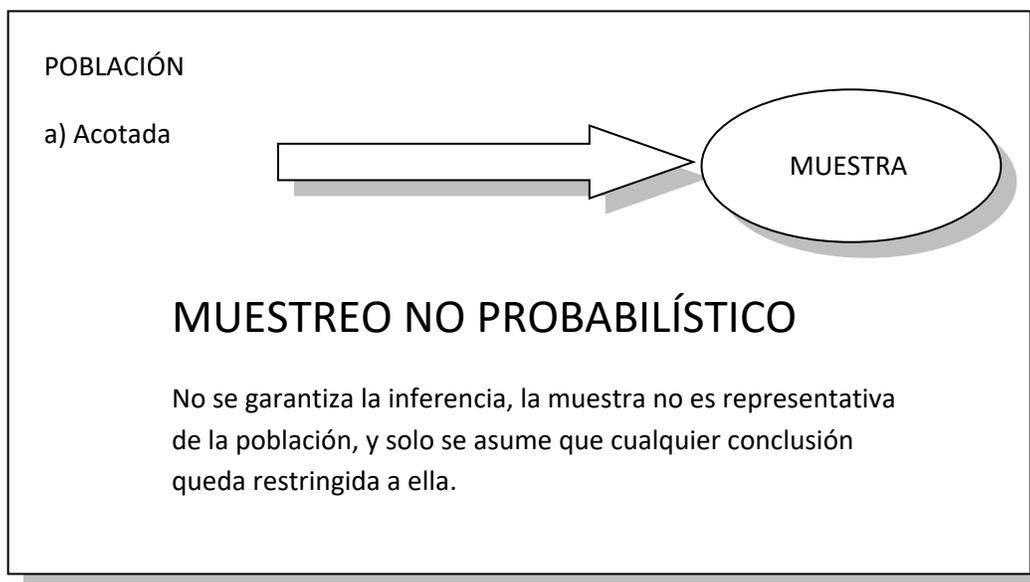
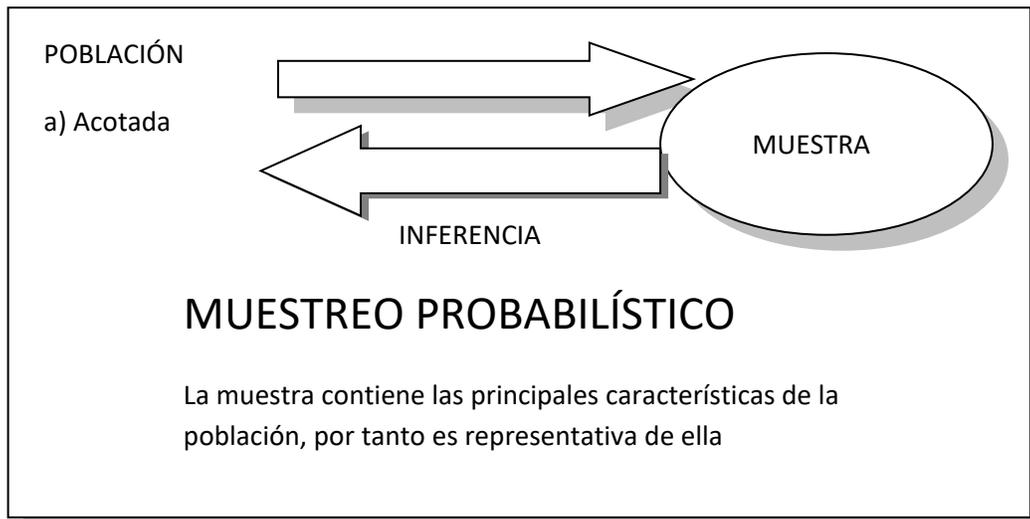
La inferencia será siempre inexacta, dado que tomamos una parte (muestra), para conocer el todo (población), aún así es posible conocer el grado de incertidumbre con el que se trabaja, y además es factible fijar el nivel de error durante el proceso en términos probabilísticos. Para que esto sea posible, la definición de población debe ser precisa para que no se produzcan sesgos desde la propia definición. Luego, si la muestra es representativa de la población, contendrá las mismas características que ella. Es conveniente aclarar en este punto, que la representatividad de la muestra estará dada en buena medida por la técnica de muestreo utilizada.

Es importante hacer notar que toda vez que calculemos un indicador sobre la población, lo denominaremos **parámetro**. Puesto que hay poblaciones que son muy difíciles de abarcar en su totalidad, utilizamos el muestreo para estimar ese parámetro; en tal caso el indicador será un **estadístico**. La estadística inferencial permite la estimación confiable de parámetros toda vez que se obtienen muestras representativas de una población; y estas solo se obtienen cuando el proceso de muestreo es de tipo probabilístico.

Sin embargo, no siempre se busca obtener una muestra representativa de la población para realizar una inferencia. En muchos casos, interesa tener una aproximación a la población pero no es finalidad de la investigación la estimación de un parámetro. Por ejemplo, en las etapas previas a la preparación de un instrumento de recolección de datos, éstos se prueban sobre muestras pequeñas de la población que no son representativas de la misma. Lo que se pretende es simplemente comprobar la confiabilidad del instrumento. Otras veces, se desea recolectar una

muestra intencionalmente construida y sacar conclusiones sobre ella, sin extrapolar el resultado a la población. Un ejemplo de esto sería el trabajo que se realiza en unidades carcelarias, donde las conclusiones del estudio quedan confinadas a la muestra estudiada. Finalmente, existen situaciones donde es muy complejo y costoso obtener una muestra representativa de la población y es necesario trabajar sobre los casos que se disponen. Por ejemplo, los estudios que se enfocan sobre características particulares que solo están presentes en una fracción muy pequeña de la población, tal el caso de enfermedades poco comunes como las encefalopatías virales y otros padecimientos similares. En situaciones donde no se pretende realizar una inferencia o cuando no es posible hacerlo, el muestreo es de tipo no probabilístico.

El siguiente esquema resume los principales aspectos del muestreo y sus propiedades.



Muestreo Probabilístico

Una cuestión fundamental para comprender el muestreo es que éste se planea y se ejecuta siempre en el contexto de un estudio con objetivos definidos. Tales estudios pueden ser de índole científica o de mercado, pero en todos los casos la finalidad del mismo determina el tipo de muestra requerida. En este apartado, realizaremos una breve descripción de los tipos de muestro conocidos como probabilísticos, es decir aquellos que son recomendados si se quiere garantizar una apropiada inferencia

estadística.

Muestreo aleatorio simple

En este tipo de muestreo, cada uno de los elementos de la población tiene la misma probabilidad de estar en la muestra seleccionada. Supongamos que definimos como población a todos los alumnos que cursan la carrera de enfermería en la Universidad Nacional de Córdoba. Mediante el registro de matrícula, sabemos que se encuentran anotados 875 estudiantes. Para denotar el tamaño de la población utilizaremos la letra **N**, en mayúscula, para diferenciarlo de la muestra, para la cual utilizaremos la letra **n**, en minúscula. Así el tamaño de la población de interés es $N=875$. De esta población, queremos tomar una parte, una muestra, que sea representativa de la misma; para ello realizamos un muestreo aleatorio simple. Nótese que bajo estas circunstancias es sencillo realizar este tipo de muestreo porque las unidades de muestreo se encuentran todas listadas de antemano (registro de alumnos matriculados). Ahora bien, para obtener de esta población una muestra, debemos definir el tamaño de la misma. Así, tendríamos que:

$N=875$ (población)

$n=100$ (muestra)

En este ejemplo, tamaño muestral lo hemos fijado de antemano, pero en una investigación real éste deberá calcularse según dos criterios que son: el tamaño del efecto y la potencia estadística. Estos criterios exceden la explicación que podemos dar en este contexto, pero cabe destacar que en las publicaciones científicas, existe casi siempre una referencia a ellos. Además, en los programas estadísticos para poder estimar el tamaño muestral necesitamos establecer una serie de parámetros tales como el intervalo de confianza de la estimación, la potencia estadística, etc. Estos temas exceden lo que podemos tratar en este tutorial, pero son fundamentales en la estadística inferencial.

Continuando con el ejemplo, supongamos que el estudio que se lleva a cabo indaga acerca de las perspectivas laborales que tienen los estudiantes. Para tal fin se aplica un cuestionario que contiene la siguiente pregunta: *“De acuerdo a su opinión, la profesión de enfermero tiene una demanda laboral: a) elevada, b) moderada, c) escasa, d) casi nula”*. Como ya se dijo, para obtener una aproximación a la opinión general de los estudiantes, es necesario que la muestra sea representativa del total de la población, y para ello se opta por un muestreo aleatorio simple. Una condición de este tipo de muestreo es que cada uno de los miembros de la población tenga la misma chance de ser seleccionado, por lo tanto para obtener la muestra de 100 estudiantes

procederíamos de la siguiente manera:

- a) numeramos a los estudiantes desde el 1 al 875,
- b) seleccionamos uno al azar, ese será el primer estudiante de la muestra,
- c) reponemos el estudiante seleccionado a la población,
- d) seleccionamos un segundo estudiante para la muestra,
- e) el procedimiento se repite hasta completar los 100 estudiantes.

Es necesario reponer el estudiante seleccionado a la población porque si así no se hiciera, el primer estudiante tiene una probabilidad de ser seleccionado de $1/875$, pero el estudiante número cien tiene una probabilidad de $1/775$, que es diferente a la del primero. El muestreo aleatorio simple, implica siempre reponer la unidad seleccionada a la población si queremos que cada individuo seleccionado tenga la misma probabilidad de ser escogido, aun a riesgo de volverlo a seleccionar. Se desprende de lo dicho que este tipo de muestreo es aplicable cuando:

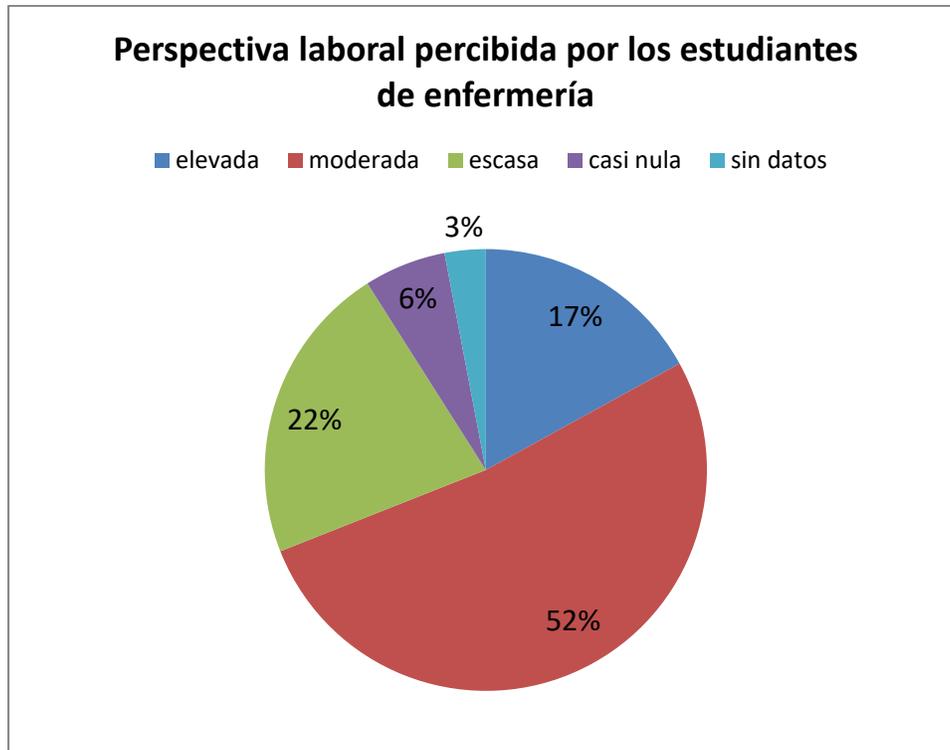
- a) la población no es muy grande,
- b) se tiene acceso a todos los individuos,
- c) es posible hacer un listado exhaustivo de la población,
- d) se puede repetir el proceso de muestreo cuantas veces sea necesario.

Supongamos entonces que la encuesta arrojó los siguientes resultados:

Perspectiva laboral percibida en estudiantes de la carrera de enfermería

elevada	moderada	escasa	casi nula	sin datos
17	52	22	6	3

Los mismos datos pueden representarse en porcentajes en un gráfico de sectores, tal como se muestra a continuación.



Podemos concluir que la mayoría (52%) estima que su título de grado le garantiza una inserción laboral moderada, en un porcentaje menor (22%) se encuentran aquellos que tiene una opinión poco optimista de tal inserción, aunque los que consideran que las oportunidades de trabajo son elevadas se aproximan a aquello que consideran nula esta posibilidad (17%). Dado que el procedimiento para recoger esta muestra ha sido aleatorio simple, es posible afirmar que de haber encuestado a todos los estudiantes, los resultados estarían muy próximos a los reportados. El sustento de tal inferencia es posible porque la muestra es representativa del total de estudiantes de la carrera de enfermería. Una conclusión inicial a la que lleva estos datos, es que casi tres cuartos de los estudiantes, estima que sus posibilidades laborales no estarían garantizadas con la obtención del título.

Muestreo aleatorio sistemático

El muestreo aleatorio sistemático es una variante del muestreo aleatorio simple para ser aplicado a poblaciones mucho más grandes. En este caso los elementos de la muestra se eligen a intervalos del listado completo de la población. Supongamos que la Cámara de Comercio, desea realizar un sondeo de opinión sobre el costo de la carga impositiva entre pequeños comerciantes. Según su registro, existen actualmente 6600 comerciantes en esa categoría; entre ellos se desea recolectar una muestra de 450

comerciantes, a quienes se le solicitará su opinión para la siguiente pregunta: *“En qué medida los impuestos que paga por su actividad comercial, restringe las posibilidades de expansión de su negocio”*; las opciones son: a) completamente, b) severamente, c) moderadamente, d) levemente.

Dado que el tamaño de la población es grande, se opta por realizar un muestreo de tipo aleatorio sistemático. Para ello, se debe escoger el valor de arranque para la selección de los participantes. Este valor se obtiene del cociente entre el total de la población y el tamaño de la muestra; esto es:

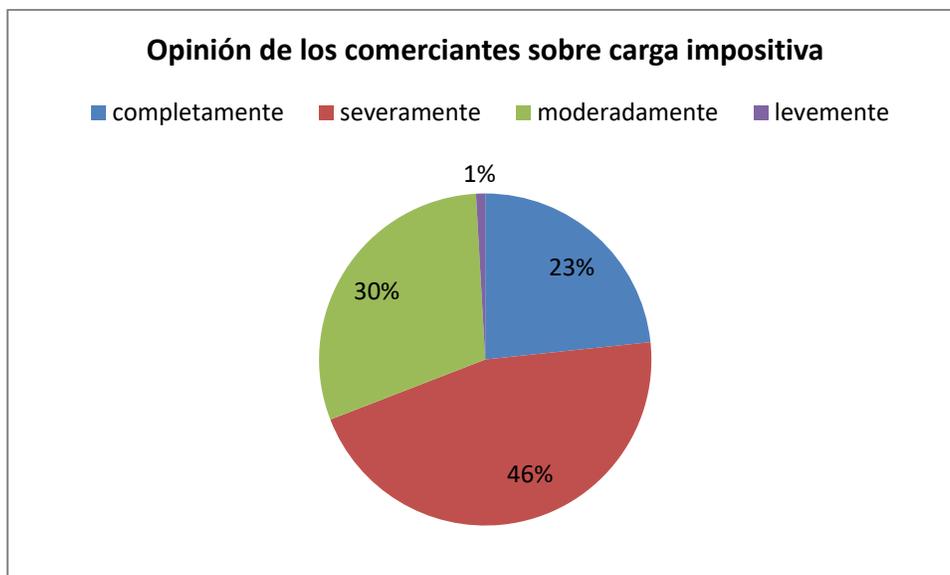
$$N/n = 6600/450 = 14,66$$

Dado que 14,66 es un número decimal lo redondeamos al entero más próximo, esto es 15. Esto nos dice que, al dividir la población sobre la muestra, sabemos que aquella puede ser parcelada 450 veces si tomamos grupos de aproximadamente 15 personas. El muestreo sistemático se basa en la elección aleatoria de un comerciante cada 15, para esto escogemos un valor por sorteo comprendido entre 1 y 15. Suponiendo que el valor sorteado es el número 7, la primera unidad de muestreo será el séptimo comerciante de la lista. El siguiente comerciante seleccionado será el que caiga en el lugar $7+15=22$, esto es, el comerciante vigesimosegundo de la lista es la siguiente unidad de muestreo; se deduce entonces que la tercera unidad de muestreo es el comerciante en el puesto $22+15=37$. Mediante este procedimiento se escogen todas las unidades de muestreo a intervalos regulares del listado original. Este procedimiento requiere el cuidado de mezclar los elementos de la lista de comerciantes (población), para que no contenga ningún sesgo. Nótese que el muestreo sistemático puede aplicarse a poblaciones más grandes al presentar menos complicaciones que un muestreo aleatorio simple; sin embargo, si este último es aplicable es el método que debe escogerse.

Continuando con el ejemplo que esbozamos, el resultado de la encuesta puede resumirse en la siguiente tabla y gráfico:

Opinión de los comerciantes sobre la carga impositiva

completamente	severamente	moderadamente	levemente
105	206	135	4



De acuerdo a los resultados obtenidos en la muestra seleccionada, es posible afirmar que la carga impositiva representa un peso importante de la actividad comercial, dado que una amplia mayoría de los comerciantes (69%) reparten su opinión en las categorías completamente y moderadamente.

Muestreo aleatorio estratificado

Como su nombre lo indica, en este tipo de muestreo es necesario definir los estratos de los cuales se van a obtener las muestras. En términos sencillos, un estrato se determina a partir de cualquier propiedad de la población que permita dividirla, tal división posibilita realizar un muestreo aleatorio dentro de cada estrato, dado que siempre serán más pequeños que la población.

Pareciera relativamente sencillo determinar qué puede definirse como estrato para una población, pero no es tarea fácil dado que cada estrato debe ser homogéneo en su interior, y además deben ser heterogéneos entre ellos. Homogeneidad y heterogeneidad deben entenderse en el contexto de la variable que se está estudiando. Volvamos por un momento al ejemplo dado anteriormente sobre la opinión de pequeños comerciantes con respecto a la carga impositiva. Supongamos que es posible establecer una separación precisa entre los comercios, de modo que estos puedan ser clasificados como grandes, medianos y pequeños. El tamaño del comercio en este caso es el estrato y por tanto todos aquellos comercios clasificados como pequeños, comparten similitudes. Siguiendo este razonamiento, podemos suponer que en otro estrato, digamos el de grandes comerciantes, su opinión será diferente respecto de los comerciantes pequeños y se evidenciará heterogeneidad en la opinión a través de los

estratos estudiados. Finalmente, es de mencionar que los estratos, una vez definidos, deben ser exhaustivos para la población.

El tipo de gestión de una escuela es un atributo que puede funcionar como estrato para la población, así es posible realizar un listado completo de los alumnos que asisten a la escuela media en una ciudad dada, considerando como estratos a las escuelas de gestión pública y de gestión privada. Si se está considerando el alumnado, las categorías descritas funcionan adecuadamente como estratos, dado que un alumno puede asistir sólo a una escuela. Tomemos por caso una variable ficticia que definimos como consumos culturales recreativos (v.g. música, lecturas, etc.). Podemos suponer en este ejemplo que los consumos culturales recreativos de los alumnos que asisten a escuelas de gestión pública serán distintos de aquellos que asisten a escuelas de gestión privada. Los estratos definidos permitirían captar el grado de heterogeneidad entre ellos. Asimismo, la gestión de la escuela permite definir un grupo más o menos homogéneo en esa variable, cuando la analizamos al interior del estrato.

Dado que en general los estratos contienen diferente número de unidades de muestreo, en ocasiones el muestreo estratificado se toma proporcional al tamaño del estrato. En nuestro ejemplo, seguramente habrá menos alumnos matriculados en escuelas privadas que en escuelas públicas, por lo tanto será necesario tener en la muestra más alumnos de este último tipo de escuelas. Sin embargo, en ocasiones no es necesario respetar la proporcionalidad de los estratos, de modo que una vez definidos los mismos, se pueden extraer muestras de igual tamaño de cada uno de ellos. La necesidad de tener muestras proporcionales al estrato, dependerá del tipo de variable que se analice. En general, cuanto menor correlación tiene la variable de interés con el estrato, menos importante es mantener la proporcionalidad del muestreo.

Veamos ahora un ejemplo de muestreo estratificado: el Organismo Consultivo Nacional para la Educación Superior, lanzó una encuesta de opinión sobre las condiciones laborales en las distintas unidades académicas del país. Se contaron 3042 docentes con dedicación exclusiva, de los cuales se decidió tomar una muestra de 600 docentes. Para realizar el muestreo se dividió la población en los siguientes estratos de acuerdo a la orientación del título dado por la unidad académica o facultad:

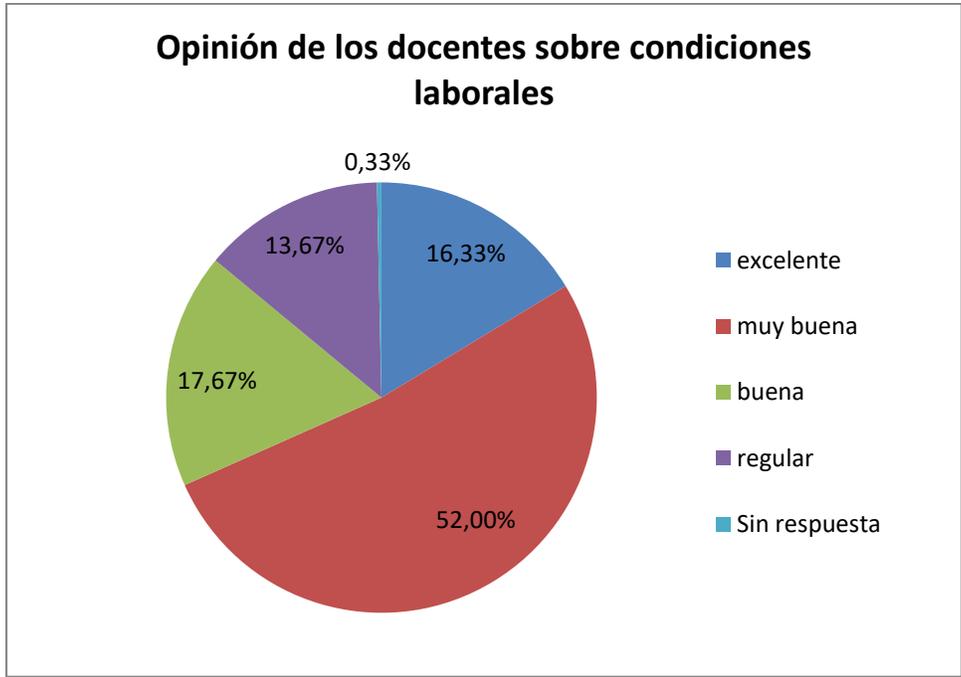
Ingeniería	Ciencias de la Salud	Ciencias Sociales	Artes y Humanidades	Matemática Física y Astronomía
660	677	578	758	369
21.7%	22.25%	19%	24.9%	12.13%
130	134	114	149	73

La segunda fila de la tabla contiene la cantidad de docentes por estrato y la tercera fila representa el porcentaje del total de la población. Si se quisiera realizar un muestreo proporcional al estrato, utilizamos ese porcentaje para conocer cuántos docentes (por estrato) deberían seleccionarse para la muestra (cuarta fila). Si la muestra requerida es de 600 docentes, multiplicamos: Proporción de docentes en el estrato, por, tamaño muestral, y así tenemos la cantidad de docentes que se requieren para ese estrato. Así, en ingeniería la cantidad de docentes a ser encuestados serían: $0,217 \times 600 = 130,2 \approx 130$; esto es, ciento treinta docentes. Este procedimiento se repite para cada estrato y de este modo, obtenemos la muestra necesaria, proporcionada al tamaño del estrato. El procedimiento nos da una aproximación al número de unidades de muestreo requeridas, pero existen otros modos de ponderación más exactos. Nótese que la multiplicación se realiza empleando la frecuencia relativa y no el porcentaje (más adelante veremos que el porcentaje se deriva directamente de la frecuencia relativa).

Continuando con el ejemplo, diremos que en el sondeo de opinión se realiza la siguiente pregunta: *“Considera Ud. que las condiciones laborales de su unidad académica son: a) excelentes, b) muy buenas, c) buenas, d) regulares”*. Las categorías de respuestas se resumen en la siguiente tabla y gráfico:

Opinión sobre las condiciones laborales de los docentes de educación superior

excelente	muy buena	buena	regular	Sin respuesta
98	312	106	82	2



Se observa que poco más de la mitad de los docentes encuestados considera que las condiciones laborales de su unidad académica son muy buenas. El resto de las opiniones muestran porcentajes similares, pero en categorías disímiles en su contenido. Dado que en este tipo de muestreo, cada estrato puede estudiarse como una unidad en sí misma, puede desagregarse la información por unidad académica, es decir en los diferentes estratos, lo cual permitiría analizarlos y compararlos.

En la tabla que se muestra a continuación se ha realizado esto, favoreciendo la comparación pretendida.

	excelente	muy buena	buena	regular	sin respuesta
Ingeniería	40	71	7	11	0
Ciencias de la Salud	18	79	27	19	1
Ciencias Sociales	12	59	17	15	0
Artes y Humanidades	11	88	45	34	1
Matemática, Física y Astronomía	17	15	10	3	0

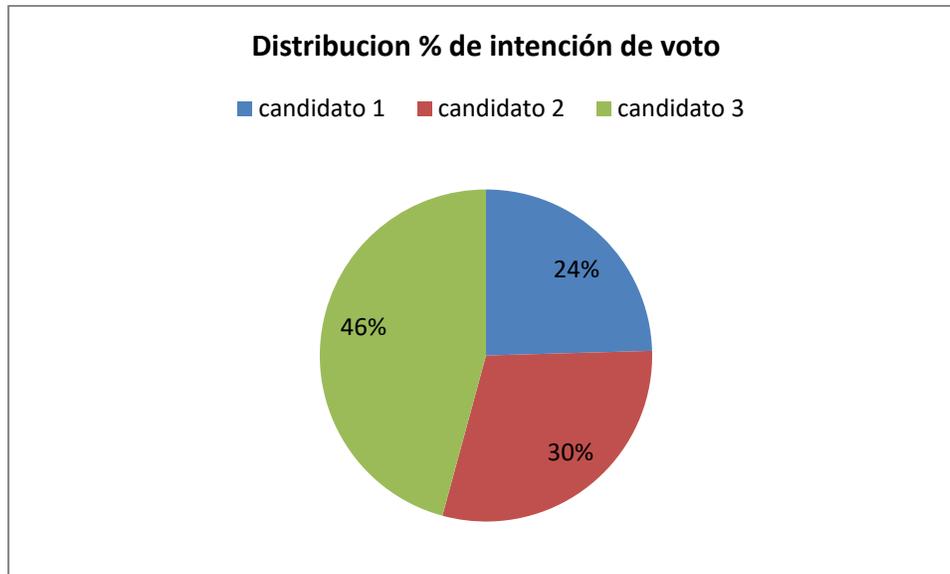
Observando cada estrato vemos que hay diferencias de opinión de los docentes de las distintas unidades académicas: por ejemplo, los docentes de la facultad de Matemática, Astronomía y Física tienden a una opinión positiva en relación a sus condiciones laborales (la mayoría de las respuestas corresponden a las categorías excelente y muy buena) mientras que más de la mitad de los docentes de Artes y Humanidades tienen una opinión positiva y el resto de las categorías más representadas son las correspondiente a buena y regular. Otra conclusión en este

mismo sentido sería que, los docentes de las Facultades de Ingeniería y Matemática, Astronomía y Física, tienen opiniones más parecidas entre sí; lo mismo que las facultades de Artes y Humanidades y Ciencias de la Salud.

Muestreo por conglomerados

En el caso en que no pueda definirse el límite de la población, o bien ésta sea demasiado grande como para obtener un listado de las unidades de muestreo, es posible trazar conglomerados dentro de los cuales realizar el muestreo. El término conglomerado se ha mantenido en la literatura puesto que proviene de la planificación para obtener grandes muestras de la población y se aplica a las parcelas geográficas o demográficas en que ésta puede dividirse. De acuerdo a la variable de interés, se intenta que los conglomerados definidos sean, en general, heterogéneos en su interior ya que deben conservar la diversidad que existe en la población, y homogéneos entre sí ya que unos pocos conglomerados deberían alcanzar para representar toda la población.

El muestreo por conglomerados se aplica cuando las unidades de muestreo no están listadas; por tanto es necesario trasladarse hasta donde se encuentran estas unidades. Como se mencionó, el conglomerado suele coincidir con una parcela geográfica. Como ejemplo, puede tomarse los diferentes distritos electorales en los sondeos de opinión. La zona de cobertura de una escuela también puede ser tomada como un conglomerado. Para comprender mejor este procedimiento, tomemos otro ejemplo ficticio. Una encuestadora realiza en período preelectoral un sondeo de opinión de los tres candidatos principales a intendente de la ciudad. Para este estudio, en el ejido municipal se identificaron cinco conglomerados. Se realiza un muestreo aleatorio de viviendas dentro de cada uno y se recogió información sobre la intención de voto de los candidatos. El resultado se expresa en el siguiente gráfico:



Como se aprecia, el candidato 3 es el que tiene una clara ventaja en la intención de voto, mientras que el candidato 1 y 2 están más próximos entre sí. La información podría completarse analizando la composición del voto por conglomerado; nótese que este procedimiento es similar al realizado con los estratos. Continuando con el ejemplo, supongamos que la composición de la intención de voto por candidato, queda expresada en la siguiente tabla:

Distribución de la intención de voto

	Conglomerado A	Conglomerado B	Conglomerado C	Conglomerado D	Conglomerado E	Total	distribución %
Candidato 1	84	322	86	90	400	982	24,55
Candidato 2	400	463	89	5	230	1187	29,675
Candidato 3	316	15	625	705	170	1831	45,775
Total	800	800	800	800	800	4000	

La tabla contiene valiosa información para comprender e interpretar los porcentajes, entre la que se cuenta la cantidad de encuestas realizadas, y el número total de respuestas positivas para cada candidato. Si analizamos la información comparando los distintos conglomerados, podemos observar que el candidato 3 tiene mayor intención de voto en la tabla general, pero no así en los conglomerados B y E, donde los otros candidatos tienen mayor aceptación. Esta observación puede hacerse extensiva para analizar la composición del voto por conglomerado de los otros candidatos. Pero, a diferencia del muestreo estratificado, aquí suponemos que los conglomerados son parecidos entre sí, por lo tanto, al analizar la composición del voto del candidato 3, sería erróneo suponer que los conglomerados B y E son distintos al resto.

Muestreo No Probabilístico

En ocasiones no es necesario realizar un muestreo aleatorio o probabilístico, dado que no es el objetivo de la investigación la inferencia a la población. Otras veces, se puede usar este tipo de muestreo en estudios orientados a una parcela específica de una población bien definida. En cualquier caso, el resultado del muestreo no se orienta a obtener datos generalizables.

Muestreo accidental

En el muestreo accidental se procede a tomar para la muestra aquellos individuos que desean participar del estudio. Un ejemplo son las encuestas callejeras, donde los encuestadores solicitan respuestas a transeúntes que prestan su colaboración. Los muestreos accidentales tienen muy poca validez y raramente se desprenden generalizaciones de ellos, aunque son excelentes para la prueba de las herramientas de testeo, tales como los cuestionarios de opinión. En la actualidad, es muy frecuente ver en los medios de comunicación este tipo de encuestas, y la información es utilizada para valorar la opinión del consumidor. Por ejemplo, en los periódicos por internet, se suele pedir al lector que valore el interés por una nota periodística, en otros sitios que ofrecen ayuda en línea, se pide que se valoren en qué medida le fue útil la ayuda ofrecida, etc. Estos son sondeos basados en muestras accidentales, que no tienen ninguna posibilidad de generalización.

Muestreo por cuotas

La metodología del muestreo por cuotas es similar a la que se utiliza en el muestreo accidental, pero en este caso se debe cumplir con una cuota en las unidades de muestreo. El término cuota se aplica a una característica que define claramente a la población que se desea muestrear. Por ejemplo, si estamos interesados en los alumnos de una escuela, éstos pueden dividirse entre varones y mujeres; si además sabemos que la proporción es 3/1 para mujeres, en el muestreo por cuotas deberemos recolectar una muestra que tenga una mayor proporción de mujeres que de hombres, definiendo la cuota mayor para las mujeres. En este ejemplo, si se desea obtener una muestra de 40 alumnos, 30 de ellos deberán ser mujeres, que sería la cuota para el sexo femenino. Como se dijo, el procedimiento del muestreo es accidental, de modo que participarán las primeras 30 mujeres, y los primeros 10 varones que den su consentimiento. Ahora bien, por tratarse de un muestreo no probabilístico, puede que no se quiera obtener un muestreo proporcional, por lo cual puede fijarse una cuota similar para los alumnos de ambos sexos. En tal caso, participarán 20 mujeres y 20 varones en la muestra.

Muestreo intencional (por conveniencia u orientado)

En este tipo de muestreo interesa la unidad de muestreo, más que el proceso por el cual se obtiene la muestra. Por ejemplo, es posible que estemos interesados en recoger una muestra de directores de escuelas primarias con más de quince años de trayectoria, o mujeres que ejerzan en cargos ejecutivos. En cualquiera de los casos, el principal interés reside en la información que aporta la unidad de muestreo en un conjunto de variables, luego puede fijarse la cantidad de unidades de muestreo que quieran seleccionarse en función de su disponibilidad.

Bola de nieve (o en cascada)

El muestreo comienza localizando a unos pocos individuos, y mediante estos se intenta localizar otros que permitan completar el tamaño de la muestra establecida. Este tipo de muestreo se emplea muy frecuentemente cuando se hacen estudios con poblaciones a las que resulta difícil acceder o que son muy poco comunes. Por ejemplo, aquellos investigadores interesados en sectas religiosas, encontrarán muy complicado realizar un muestreo probabilístico, pero puede facilitar su tarea de campo si logra contactar un miembro de la secta, y a través de él, llegar a otros. El estudio de los sociolectos es otro ejemplo donde se utiliza este tipo de muestreo. Primero se contacta a un miembro del grupo de interés y a través de él se arma una red de contactos que sean potenciales integrantes de la muestra. La conformación de grupos focales es otro ejemplo en donde utiliza este muestreo.

Comentarios Finales

En los apartados sobre la metodología empleada para obtener una muestra probabilística de la población, suele presentarse el error de muestreo como un estadístico. Éste término no debe confundirse con un muestreo erróneo, sino que se refiere al grado de precisión en que la muestra representa la población. El error de muestreo suele expresarse en porcentaje, y se espera que este no supere el 5%. En tal caso, estaríamos seguros en un 95% de que el estadístico calculado en la muestra, coincide con el parámetro poblacional. El término proviene de la siguiente comprobación empírica: si medimos una variable y promediamos su valor en varias muestras iguales de una misma población, extraídas mediante un muestreo probabilístico; la distribución de los promedios sigue un modelo normal. Tomando como base ese modelo, es factible aproximarse con bastante fiabilidad al nivel de error cometido en la muestra obtenida. El modelo de distribución normal se verá más

adelante.

Vale decir entonces que el procedimiento de muestreo probabilístico garantiza que ese error sea el menor posible, y sólo en este tipo de muestreos es posible calcularlo. Por otra parte, los errores en la recogida de una muestra se denominan sesgos muestrales. Existen principalmente dos razones por las cuales se producen sesgos en la muestra, una de las cuales es haber definido mal la población hacia la cual se quiere dirigir la inferencia. Recuérdese que existen poblaciones acotadas y no acotadas, por lo tanto una mala definición de la población puede repercutir en que se trabaje sobre límites ficticios de la misma. Por ejemplo, si se define como población de interés aquellas personas que han sufrido asalto a mano armada en su vivienda y se toma como base las denuncias realizadas a la policía, deberá tenerse en cuenta que dichas denuncias no contienen el total de la población, dado que varios de esos hechos delictivos no se denuncian. Además, puede que la variable de interés quede moderada por otra variable que no se ha tenido en cuenta, por ejemplo, que las denuncias de esos hechos delictivos sean más frecuentes en ciertas zonas de la ciudad. Pretender una generalización fiable bajo estas circunstancias es riesgoso, puesto que estaríamos cometiendo un error no reconocido en la inferencia.

Otra razón que conduce a sesgos en las muestras, afectan a la definición de estrato y conglomerado. Los estratos deben ser homogéneos en su interior y heterogéneos entre ellos. Por ejemplo, estudiantes de primaria y de secundaria podrían definir dos estratos que cumplirían con ese requerimiento, dado que la diferencia de edad garantiza que tengan una composición desigual entre ellos y sean parecidos en su interior. Pero no ocurriría lo mismo si definiéramos dos estratos como alumnos de tercero y cuarto grado. Es muy probable que en conjunto, los alumnos tengan más similitudes que diferencias, y en tal caso, ello invalidaría la definición de estrato. Algo similar ocurre cuando se trazan los conglomerados, se supone que estos pueden ser diferentes en su interior y tener características similares entre ellos. Por ejemplo, si se divide el mapa de Córdoba en nueve conglomerados, siete de los cuales tienen composiciones de familias de clase baja, y los dos restantes contienen el resto de la población. Si se decide tomar una muestra sobre sólo tres conglomerados, lo más probable es que se favorezcan sectores de clase social desfavorecida. En una situación como la descrita, será necesario redefinir los límites de los conglomerados.

Misceláneas

La encuesta del Literary Digest: En 1936 se llevaron a cabo elecciones en EE.UU. y los candidatos a presidentes eran Alf Landon y Franklin D. Roosevelt. La revista *Literary Digest* llevó a cabo una encuesta y predijo que el nuevo presidente sería Alf Landon. Los resultados de la elección dieron por ganador a Franklin D. Roosevelt por una gran mayoría. ¿Qué fue lo que falló? Como luego se evidenció, la revista envió la encuesta

por correo a 10 millones de ciudadanos estadounidenses, la cual es una muestra lo suficientemente grande como para ser representativa, pero solo contestaron la encuesta 2.3 millones de votantes, lo cual representa solo el 23% de los encuestados. Esta cifra es demasiado pequeña como para avanzar alguna predicción fiable del resultado de la votación. Aún así, los directivos del Literary Digest se aventuraron a predecir un resultado electoral, en el que se daba por ganador al candidato Alf Landon. La pérdida de credibilidad de la revista y la repercusión del fracaso en las encuestas electorales, condujo a que se modificaran sustantivamente las técnicas de muestreo en los sondeos de opinión electoral. Una cuestión que quedó claro fue que: no es el tamaño de la muestra lo que garantiza su representatividad, sino el método empleado para recogerla.

La paradoja del candidato: Hay importantes cuestiones sobre las elecciones por sistemas de votación que han sido estudiadas por matemáticos y estadísticos, tales como Richard, G. Nemi, William H. Riker, Donald Saari y Allyn Jackson. Una de ellas se refiere a la paradoja del candidato la cual expondremos brevemente. Supongamos que hacemos competir en una elección presidencial al partido del CENTRO, el de DERECHA y el de IZQUIERDA; ¿es posible que de antemano arreglemos la elección para que gane uno de ellos? En teoría no sería posible porque depende del caudal de votos que pueda captar cada partido, pero analicemos esta situación hipotética. Se les pide a los votantes que ordenen a los candidatos en relación a su preferencia, y de ello resulta que:

- a) 1/3 de la población prefiere C – D – I,
- b) 1/3 de la población prefiere D – I – C, y
- c) 1/3 de la población prefiere I – C – D.

En el caso a) decimos que, un tercio de la población prefiere al candidato de centro, por sobre el candidato de derecha y a este por sobre el candidato de izquierda. Similar interpretación se aplica en los casos b) y c). Supongamos que deseamos que el candidato favorecido sea el de izquierda, entonces en primera vuelta hacemos competir al candidato de centro y de derecha. De esta elección resulta favorecido el candidato del centro con 2/3 de los votos, y al candidato de derecha lo elige solo 1/3 de la población. En segunda vuelta hacemos que compitan el candidato de centro con el de izquierda, y tenemos que 2/3 de la población prefiere al candidato de izquierda por encima del de centro, mientras que solo 1/3 de la población prefiere al candidato del centro por sobre el de izquierda. El resultado de estas dos vueltas electorales ¡consagra ganador al candidato de izquierda! Lo interesante de este juego es que es posible hacer lo mismo con los otros candidatos.

Esta paradoja es posible porque en realidad no existe una relación transitiva entre los candidatos de centro, izquierda y derecha. Las relaciones transitivas solo se verifican cuando es posible ordenar atributos o personas con cuantificadores tales como mayor que/menor que o similares. En tal caso, si decimos que $A > B$ (A es mayor que B) y luego $B > C$ (B mayor que C), entonces se deduce que A es mayor que C por propiedad transitiva. De ello se desprende que esta propiedad no puede aplicarse cuando la relación entre atributos es del tipo, "prefiere a", dado que este no es un cuantificador. Es decir, si solo se trata de ordenar por preferencia, se puede dar que alguien prefiera A sobre B, B sobre C, pero a C sobre A. La paradoja en este caso se conoce como Paradoja de Arrow, en honor a Kenneth J. Arrow, quien postulo que un sistema de votación democrática perfecto es imposible.

Capítulo 3

Tablas de Frecuencia

Las tablas de frecuencia representan una manera sencilla de resumir información sobre una (o varias) variable y se basan en el conteo de la cantidad de unidades que quedan comprendidas en los diferentes niveles de medición de la misma. Las tablas casi siempre acompañan a los textos y son las referencias que se toman cuando es necesario mostrar comparaciones o destacar algún fenómeno puntual. Una tabla de frecuencia se construye a partir de filas y columnas, y por conveniencia, se utiliza la primera para los valores de la variable y las segundas para los distintos tipos de frecuencia que es posible calcular. En los apartados que siguen se mostrarán ejemplos de tablas de frecuencia para los sistemas de medición nominal, ordinal y métrico.

Tabla de frecuencia: sistema de medición nominal

Supongamos que un investigador ha realizado una medición de la variable estado civil, en el conjunto de todos los empleados de una fábrica textil, donde $N=345$. Los niveles de medición de la variable estado civil son: a) Solteros, b) Casados, c) Separados, d) Viudos, e) Otras. Para construir una tabla de frecuencia, es necesario determinar por conteo cuántas personas se encuentran comprendidas en cada una de las categorías; la tabla de frecuencia para este ejemplo podría tener la siguiente distribución de casos:

Estado civil del personal de planta (n=345)

Variable Estado Civil	Frecuencia Absoluta
Solteros	122
Casados	107
Separados	110
Viudos	1
Otras	5
Total	345

Nótese que en la columna frecuencia absoluta realizamos el conteo de casos en cada una de las categorías de la variable, y así establecemos que en la planta existen 122 personas solteras, 107 de ellas están casadas, etc. Puesto que el procedimiento consiste en contar unidades (personas en este caso), en cada uno de los niveles de la variable, la suma de la columna frecuencia absoluta debe ser igual al total de casos.

Las tablas de frecuencia también pueden contener las frecuencias relativas y estas se pueden transformar a porcentajes. Trabajar con porcentajes es más operativo que hacerlo con las frecuencias absolutas, dado que éstos nos permiten realizar comparaciones directas en caso de tratar con otra tabla que contenga información sobre la misma variable, pero con distinto número de casos. La frecuencia relativa se obtiene del cociente entre la cantidad de casos en cada uno de los niveles de la variable, sobre el total de casos. El porcentaje, se obtiene de multiplicar la frecuencia relativa por 100. En la siguiente tabla se grafica el procedimiento, y luego se presenta resumido:

Estado civil del personal de planta (n=345)

Variable Estado Civil	Frecuencia Relativa	Porcentaje
Solteros	$122/345=0.3536$	$0.3536 \times 100 = 35.36\%$
Casados	$107/345=0.3101$	$0.3101 \times 100 = 31.01\%$
Separados	$110/345=0.2898$	$0.2898 \times 100 = 28.98\%$
Viudos	$1/345=0,0028$	$0,0028 = 0.28\%$
Otras	$5/345=0,0144$	$0,0144 \times 100 = 1.44\%$

Resumen

Variable Estado Civil	Frecuencia Relativa	Porcentaje
Solteros	0.3536	35.36%
Casados	0.3101	31.01%
Separados	0.2898	28.98%
Viudos	0,0028	0.28%
Otras	0,0144	1.44%
Total	$0.9707 \approx 1$	$97.07\% \approx 100\%$

Vamos a tomar los valores de la segunda tabla, donde se encuentra resumido el procedimiento. Nótese que la frecuencia relativa y el porcentaje representan la misma información, pero dado que el porcentaje es un valor más familiar podríamos concluir que el 35.36% del personal de planta está compuesto por personas solteras, que el 31.01% de ellos son casados, etc. Como se aprecia, la información absoluta puede expresarse de una manera diferente sin perder valor.

Dado que para obtener las frecuencias relativas, debemos dividir por el total de casos, siempre obtendremos un número decimal no periódico. En tales casos se deberán utilizar solo algunos decimales; por ejemplo, el valor real de la frecuencia relativa para solteros es de: 0,35362318840579710144927536231884, un número con más de treinta cifras en su valor decimal. Puesto que tanta precisión no es necesaria, se suele usar tres o cuatro decimales, tal como se hizo en la tabla. En tal caso, y si no se ha redondeado el valor decimal, la suma a través de la columna de frecuencias relativas dará un número muy próximo a 1. Puesto que el porcentaje es la

frecuencia relativa multiplicada por cien, la suma a través de la columna porcentaje se aproximará a 100.

Otra forma de trabajar con una tabla de frecuencia es a través de sus frecuencias acumuladas. Para obtener las frecuencias acumuladas, se deben sumar de manera descendente los valores de frecuencia en cada una de las celdas de la categoría de la variable. En las siguientes tablas se presenta el procedimiento y el resumen del mismo, utilizando como base la frecuencia absoluta y el porcentaje:

Estado civil del personal de planta (n=345). Frecuencia absoluta acumulada

Variable Estado Civil	Frecuencia Absoluta	Frecuencia Absoluta Acumulada
Solteros	122	122
Casados	107	122+107=229
Separados	110	229+110=339
Viudos	1	339+1=340
Otras	5	340+5=345
Total	345	

Porcentaje acumulado

Variable Estado Civil	Porcentaje	Porcentaje Acumulado
Solteros	35.36%	35.36%
Casados	31.01%	35.36+31.01=66.37%
Separados	28.98%	66.37+28.98=95.35%
Viudos	0.28%	95.35+0.28=95.63%
Otras	1.44%	95.63+1.44=97.07%
Total	97.07%\cong100%	

El procedimiento para obtener la frecuencia acumulada es el mismo, tanto si se trata de la frecuencia absoluta o el porcentaje. Véase la tabla de frecuencia absoluta; en la categoría de la variable solteros, se cuentan 122 individuos; puesto que no hay ninguna categoría por encima, en este punto solo se cuentan esos individuos. La siguiente categoría es casados, y en ella se cuentan 107 individuos, que sumados a los 122 de la categoría solteros dan un total de 229 personas; por lo tanto la frecuencia absoluta acumulada al contar solteros y casados es de 229 personas. Si a estas le sumamos la categoría separados, debemos añadir a las 229 personas contabilizadas, otras 110, que es la frecuencia para la categoría solteros, y entonces tendremos 339 personas contabilizadas en las tres primeras categorías de la variable. El procedimiento se repite para las dos categorías restantes, y en la última categoría que sumemos, obtendremos el total de casos. Con el porcentaje se ha procedido de la misma manera, y en resumen la tabla se presentaría de la siguiente forma:

Estado civil del personal de planta (n=345). Resumen

Variable Estado Civil	Frecuencia Absoluta Acumulada	Porcentaje Acumulado
Solteros	122	35.36%
Casados	229	66.37%
Separados	339	95.35%
Viudos	340	95.63%
Otras	345	97.07%

Una conclusión que puede derivarse de la observación de los porcentajes acumulados es que más de la mitad del personal de la planta son personas solteras y casadas, que junto con los individuos separados conforman la mayoría del personal.

Utilidad de la tabla de frecuencia a través de un ejemplo

En la empresa textil se desea adaptar la cobertura social de sus trabajadores, permitiendo la opción de contratar el servicio en distintas prestadoras de salud y de acuerdo a las reales necesidades del trabajador. Para cumplir con este objetivo, le solicita a un investigador social, que sugiera algunas recomendaciones a los efectos de que la propuesta pueda ser bien recibida por los empleados. Dadas las circunstancias, el investigador supone que la principal preocupación del asalariado en relación a su cobertura social, sería la prestación extendida al grupo familiar. Por ende en una primera etapa, realiza un relevamiento del estado civil de los 345 trabajadores. Los resultados obtenidos se presentan en la siguiente tabla, de la cual se pueden sacar las conclusiones que siguen:

Resumen

Estado Civil	Frecuencia Absoluta	Porcentaje	Porcentaje Acumulado
Solteros	122	35.36%	35.36%
Casados	107	31.01%	66.37%
Separados	110	28.98%	95.35%
Viudos	1	0.28%	95.63%
Otras	5	1.44%	97.07%
Total	345	≈100%	

- a) De los planes de salud al menos el 31% de ellos, debería incluir una ampliación para el cónyuge.
- b) Aproximadamente el 60% de los planes de salud deberá incluir en el grupo familiar a uno o más hijos. Contando con que el 31% ya incluya a un cónyuge.
- c) Al menos el 35% de los planes de salud podría incluir a un solo individuo como beneficiario.

Tabla de frecuencia: sistema de medición ordinal

Cuando se trabaja con una variable de tipo ordinal con las categorías expresadas como etiquetas, el procedimiento para construir una tabla de frecuencia es idéntico al que se usa para una variable medida en escala nominal. Si la variable se expresa como una escala numérica, cada valor es equivalente a la etiqueta de la variable. Veamos el siguiente ejemplo: un sociólogo desea medir el grado de cohesión de un grupo y además el sentido con que identifica las interacciones entre individuos. Para ello, desarrolla dos escalas likert a partir de los siguientes reactivos:

- a) Exprese su grado de acuerdo con la siguiente afirmación: “el grupo al que pertenezco se encuentra muy unido”. Utilice según su criterio una escala de 1 a 10, siendo: 1=completamente en desacuerdo; 10= completamente de acuerdo.

- b) ¿Juzgue de que manera Usted percibe las interacciones entre los miembros del grupo? Utilice según su criterio una escala de -5 a +5, siendo:
-5= los individuos interactúan de modo completamente agresivo;
+5= los individuos interactúan de modo completamente colaborativo;

Nótese que en este ejemplo, se han utilizado dos tipos diferentes de escalas, pero siguiendo un criterio similar. En el segundo caso, el signo positivo o negativo indica el polo al que tienden las interacciones, agresivas o colaborativas. Suponiendo que el investigador planea realizar una serie de tareas con la finalidad de afianzar la pertenencia grupal, necesitará al menos dos mediciones para verificar si éste aspecto del grupo se ha modificado. Tales mediciones se tomaran en antes y después de las tareas de integración. Continuando con el ejemplo, la tabla que se muestran a continuación refleja lo sucedido:

Grado de cohesión del grupo n=55

Escala de medición	Antes		Después	
	Frecuencia Absoluta	Porcentaje	Frecuencia Absoluta	Porcentaje
1	1	1.81	0	0
2	3	5.45	2	3.63
3	2	3.63	3	5.45
4	8	14.54	5	9.09
5	14	25.45	12	21.81
6	12	21.81	11	20
7	10	18.18	19	34.54
8	2	3.63	1	1.81
9	2	3.63	1	1.81
10	1	1.81	1	1.81
Total	55	99,94	55	99,95

Vamos a proponer una posible interpretación de la anterior tabla frecuencias: en primer lugar se observa que las frecuencias tienden a concentrarse en los valores medios de la escala propuesta, la cual refleja el grado de cohesión percibido por los miembros del grupo. En este sentido, las frecuencias están en los valores 5, 6 y 7; al estar estos valores en el rango medio, es posible suponer que el grupo en general se percibe medianamente cohesionado. Ahora bien, se proponen una serie de actividades que se espera que fortalezcan la cohesión del grupo, por lo tanto, un cambio en la variable, debería quedar reflejada en una distribución de frecuencias concentrada en los valores altos de la escala. Partiendo de esta interpretación, la columna de la tabla de frecuencias que refleja lo sucedido luego de la intervención del sociólogo y sus actividades grupales, muestra que las frecuencias altas siguen concentradas en los valores medios de la escala, apareciendo como diferencia entre la situación antes y después, la acumulación de casos en el valor 7 de la escala luego de la intervención (34.5%). En general, las distribuciones de frecuencias para la situación antes y después no presentan variaciones notables. Esta distribución de datos, sirve como disparador de preguntas de investigación que podrían ser las siguientes:

- a) ¿han sido adecuadas las intervenciones planeadas?
- b) ¿el tiempo transcurrido entre las mediciones antes – después ha sido el suficiente como para apreciar un cambio en el grupo?

- c) ¿se trata de un grupo reactivo al cambio?
- d) ¿los individuos han logrado identificar adecuadamente las acciones que tienden a la cohesión?
- e) ¿el instrumento de medición es el apropiado para captar un cambio sutil pero significativo en la cohesión grupal?

Supóngase que en la pregunta sobre el modo de interacción, veinte de los individuos encuestados no dan respuesta. Bajo esta circunstancia, una escala de once categorías ordenadas como la que se propone (nótese que aquí se proponen cinco valores positivos, cinco negativos y un cero para el neutro), puede que contenga varias celdas vacías o con frecuencia cero. Bajo estas circunstancias, sería conveniente para el investigador reducir las categorías de la escala, una forma de hacerlo es agrupando en categorías los valores originales. La tabla que sigue es un posible ejemplo de cómo podría hacerse.

Categoría	Las interacciones son en su mayoría agresivas	A veces las interacciones son agresivas	Las interacciones son neutras	A veces las interacciones son colaborativas	Las interacciones son en su mayoría colaborativas
Puntaje Original	-5	-3	0	1	4
	-4	-2		2	5
		-1		3	

Al reducir de este modo la escala se evitaría trabajar con categorías con frecuencia cero. Como se aprecia, el procedimiento de reducción de la escala consiste en agrupar los valores de ésta en etiquetas que reflejen los valores originales, es decir, que respete en sentido del ordenamiento de la escala. El resultado de tal agrupamiento produce una escala con cinco categorías ordenadas, cuya regla de asignación es contar los individuos que puntuaron en la escala con -5, -4 en la categoría: Las interacciones son en su mayoría agresivas, luego los individuos que puntuaron en la escala -3, -2 y -1 se agrupan en la categoría: A veces las interacciones son agresivas; el procedimiento se repite para todos los valores observados.

Con la escala de medición reducida, es posible construir una tabla de frecuencias que resulte visualmente manejable y posibilite una aproximación más sencilla a su interpretación. En la siguiente tabla, el investigador presenta los datos totales del grupo, y en columnas donde se han separado las respuestas por el género de los encuestados:

Percepción de las interacciones del grupo n=25

Las interacciones en el grupo son:	Total	Varones	Mujeres
En su mayoría agresivas	1	1	0
A veces son agresivas	7	6	1
Son neutras	6	2	4
A veces son colaborativas	7	1	6
En su mayoría colaborativas	4	2	2
Total	25	12	13

Para interpretar la información de la tabla debemos considerar primeramente la pérdida de respuestas, dado que involucra un número importante de individuos. Deberíamos concentrarnos en determinar si la pregunta ha sido apropiada, si hay variables externas al estudio que estuvieran determinando que las personas no respondan, etc.

Una línea interpretativa, podría llevarnos a comparar la distribución de frecuencias del total, con la que aparece segmentada por género. En este caso se aprecia que las frecuencias más altas están en el rango medio de la escala, con una tendencia hacia el polo colaborativo de esta. Véase la distribución de frecuencias del total y tenemos que prácticamente la misma cantidad de individuos han percibido las interacciones grupales como: A veces son agresivas, Son neutras, A veces son colaborativas. En esas categorías se contabilizan 20 individuos; de los 5 restantes que forman el total, 4 han puntuado en la categoría En su mayoría colaborativas.

Si se proyecta esa tendencia sobre las frecuencias discriminadas por género, aparecen diferencias que merecen ser atendidas: a) la frecuencia más alta en el grupo de varones se encuentra en la categoría A veces son agresivas, b) en el grupo de mujeres las frecuencias tienden a agruparse en las categorías Son neutras, A veces son colaborativas. Algunas preguntas importantes que se desprenderían de estos datos podrían ser:

- a) ¿qué factores habrían influido en la falta de respuesta a la pregunta?
- b) ¿existe verdaderamente diferencias de género en las respuestas, o simplemente se deben a la escasa cantidad de personas que respondieron?
- c) ¿las personas que respondieron son las que más interacciones con el resto del grupo han mostrado?
- d) ¿las personas que no han respondido son quienes desplegaron menos interacciones con el resto del grupo?

Tablas de frecuencias para variables métricas

Las variables métricas utilizan un continuo de valores y tienen propiedades en las que los números admiten las operaciones matemáticas. Conviene recordar que este tipo de variables pueden ser continuas si admiten cualquier valor, o bien discretas, en cuyo caso solo se admiten números enteros. Por lo tanto, las tablas de frecuencia de estas variables suelen agruparse en intervalos dado que entre dos valores cualesquiera, pueden encontrarse otros tantos valores fraccionarios, o pueden tener un rango de valores muy amplios aún siendo discretas.

Ejemplo: en una escuela se aplicó una prueba de madurez lectora (Prueba de figuras inversas de Edfeldt), para determinar el nivel de aprestamiento de los alumnos al comienzo del ciclo lectivo en el primer grado. La prueba evalúa el nivel de madurez visoespacial del alumno para el aprendizaje de las letras del alfabeto. Se evaluó una muestra de 75 niños, y el resultado de dicha evaluación se muestra en la siguiente tabla:

Valor Real del Intervalo	Valor Aparente del Intervalo	Frecuencia	Frecuencia acumulada	Porcentaje	Porcentaje acumulado
$66,0 < x \leq 68,0$	66 - 68	5	5	6,66	6,66
$69,0 < x \leq 71,0$	69 - 71	8	13	10,66	17,32
$72,0 < x \leq 74,0$	72 - 74	10	23	13,33	30,65
$75,0 < x \leq 77,0$	75 - 77	11	34	14,66	45,31
$78,0 < x \leq 80,0$	78 - 80	15	49	20	65,31
$81,0 < x \leq 83,0$	81 - 83	19	68	25,33	90,64
$84,0 < x \leq 86,0$	84 - 86	7	75	9,33	99,97 \approx 100
		75			

La interpretación de los valores de frecuencia y porcentaje son los mismos que para variables medidas en escalas nominales y ordinales. Lo que se modifica en este caso es la manera de presentar los valores de la variable.

Supongamos que los 75 escolares evaluados hubieran obtenido un valor diferente en la prueba; si en tales circunstancias construyéramos una tabla de frecuencia con cada valor individual, ésta contendría 75 renglones con valor de frecuencia igual a 1. Una tabla de estas características sería inútil. Por ello, siempre que se trabaja con variables métricas es conveniente agrupar los valores en intervalos.

La tabla que estamos presentando contiene intervalos de valor 1, esto es a los fines didácticos pues en tablas reales los intervalos son más amplios. Los valores de la prueba Edfeldt son discretos, es decir, la puntuación de prueba no contiene decimales. Para construir una tabla de frecuencia comenzamos con la primera columna que contiene los valores reales del intervalo, allí se cuentan todas las puntuaciones con valor 66 hasta 68 inclusive. Por tratarse de una variable discreta sabemos que ese

intervalo solo contiene los siguientes valores: 66, 67, 68. Si estuviéramos trabajando con variables de tipo continuas, el intervalo podría contener cualquier valor entre 66,0... y 68,9... (volveremos sobre esto más adelante). La siguiente columna representa el valor aparente del intervalo, que no es otra cosa que el valor real sin los decimales. En este caso es sencillo ver los límites del intervalo por tratarse de una variable discreta. Las columnas que siguen contienen información que ya conocemos.

Cuando se publica una tabla de frecuencia de una variable métrica (sea continua o discreta), se emplea la columna de valor aparente del intervalo puesto que es más fácil de interpretar. Tomemos el tercer intervalo cuyo valor aparente es 72 - 74, su frecuencia es 10 lo cual indica que hay diez escolares que han obtenido un puntaje que está entre el valor 72 y el valor 74.

Veamos ahora algunas precisiones contenida en la columna valor real del intervalo; existen diez escolares que han obtenido un puntaje 74 o mayor, y menor o igual que 76, esto está expresado por el límite real de intervalo $74,0 < x \leq 76,0$, donde x indica el valor obtenido por cualquiera de los diez escolares. Cuando los valores de una variable métrica se agrupan en intervalos, no es posible identificar el verdadero valor individual. Dicho en otras palabras, al agrupar los valores de la variable en intervalos se logra reducir la información para que pueda caber en una tabla que sea fácilmente manejable, el costo de ello es que se debe sacrificar ciertas precisiones en los valores individuales. Por lo tanto, la cantidad de intervalos y su amplitud, debe ser tal que permita un óptimo manejo de los datos en la tabla.

En el siguiente cuadro se muestran algunas referencias para la prueba, los cuales sirven en general para interpretar el valor obtenido por un individuo particular. Por ejemplo, si un niño obtiene un puntaje igual o menor a 60, indica que sus aptitudes visoespaciales para la lectura aún no se encuentran maduras. En el otro extremo, un escolar que obtuviera un puntaje igual o mayor a 75, mostraría que su capacidad visoespacial se encuentra madura para iniciar el aprendizaje de la lectura.

Puntajes de prueba	Interpretación
60 o menos	El niño no posee aún las aptitudes suficiente para el aprendizaje del alfabeto
74 o menos	El niño posee aptitudes limitadas para el aprendizaje del alfabeto
75 o más	El niño posee aptitudes suficientes para el aprendizaje del alfabeto

*Las interpretaciones ofrecidas son una simplificación de los estándares de la prueba

Según muestra la tabla 10, ningún escolar obtiene un puntaje menor de 60, lo cual indica que todos poseen aptitudes para el aprendizaje de alfabeto. Luego, tenemos que 23 escolares obtienen puntajes entre 66 y 74, valor que resulta del conteo de los tres primeros intervalos de la variable, y que se verifica en la columna de frecuencia acumulada. Podemos decir, según los valores de referencia de la prueba, que un

30,65% de la muestra de escolares tiene aptitudes limitadas para el aprendizaje del alfabeto, por lo cual hay que trabajar en este grupo reforzando el aprendizaje visual de los grafemas. El 69,35% de los niños posee aptitudes suficientes para el aprendizaje del alfabeto. Con esos datos concluimos que es posible focalizar diferencialmente el proceso de enseñanza para aquellos niños que aún no han desarrollado suficientemente sus capacidades visoespaciales.

Veamos ahora un ejemplo con una variable métrica continua y las diferencias en los intervalos de confianza. En una escuela se les tomó un examen a los alumnos y se los calificó con una nota, también se contabilizó el tiempo que emplearon para realizarla. En la siguiente tabla se muestra el tiempo empleado por los alumnos agrupados en intervalos. Una aclaración para esta tabla es que el tiempo está contado en minutos.

Valor Real del Intervalo	Valor Aparente del Intervalo	Frecuencia	Frecuencia acumulada	Porcentaje	Porcentaje acumulado
$60,0 < x \leq 65,59$	60 - 65	2	2	2,78	2,78
$66,0 < x \leq 71,59$	66 - 71	12	14	16,67	19,45
$72,0 < x \leq 77,59$	72 - 77	16	30	22,22	41,67
$78,0 < x \leq 83,59$	78 - 83	21	51	29,17	70,84
$84,0 < x \leq 89,59$	84 - 89	14	65	19,44	90,28
$90,0 < x \leq 95,59$	90 - 95	7	72	9,72	100
		72		100	

La interpretación que podemos ofrecer de estos datos sigue la misma lógica que la aplicada al ejemplo anterior. Aquí nos detendremos en la diferencia entre valor real de intervalo y valor aparente. Si tomamos el valor del primer intervalo vemos que dos individuos emplearon entre 60 y poco más de 66 minutos en realizar la prueba. El valor exacto del tiempo de prueba de estos dos individuos lo desconocemos, pues es información que se pierde al agrupar en intervalos. Luego sabemos que 12 escolares tardaron entre 66 y poco más de 71 minutos en finalizar la prueba.

Los intervalos primero y segundo en su valor real se representan como $60,0 < x \leq 65,59$ y $66,0 < x \leq 71,59$ lo que indica que el primero abarca los tiempos de 60 minutos hasta 65 minutos, 59 segundos. Dado que la tabla se construyó con el tiempo en minutos, transcurridos los 65 minutos con 59 segundos, pasamos al siguiente intervalo que comienza en 66 minutos y culmina en 71 minutos y 59 segundos. Así continuamos para cada uno de los intervalos. Este modo de presentar la información es muy preciso pero engorroso, por lo cual simplificamos mucho la escritura de una conclusión si solo empleamos el valor aparente del intervalo, que sería 60 - 65 y 66 - 71 para el primer y segundo intervalo respectivamente. De este modo, podríamos redactar una conclusión diciendo que dos alumnos fueron los que emplearon el menor tiempo para realizar la prueba, que estuvo entre 60 y 65 minutos. Quienes más demoraron en realizarla (siete alumnos) tardaron entre 90 y 95 minutos. La mayoría

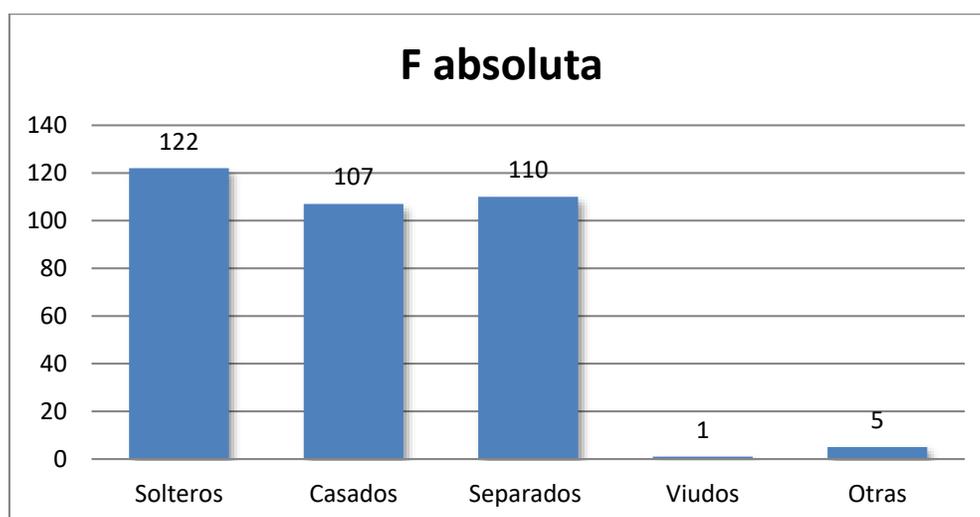
empleó un tiempo de 78 a 83 minutos (21 alumnos).

Representación gráfica de las tablas de frecuencias

En las páginas que siguen se mostrarán algunos de los gráficos más usados para representar los datos contenidos en las tablas de frecuencias. A modo de resumen, se usaran los mismos ejemplos dados anteriormente.

Diagrama de barras

El siguiente diagrama representa la primera columna de la tabla que pertenece a la frecuencia absoluta de la variable estado civil de los empleados de una planta textil.



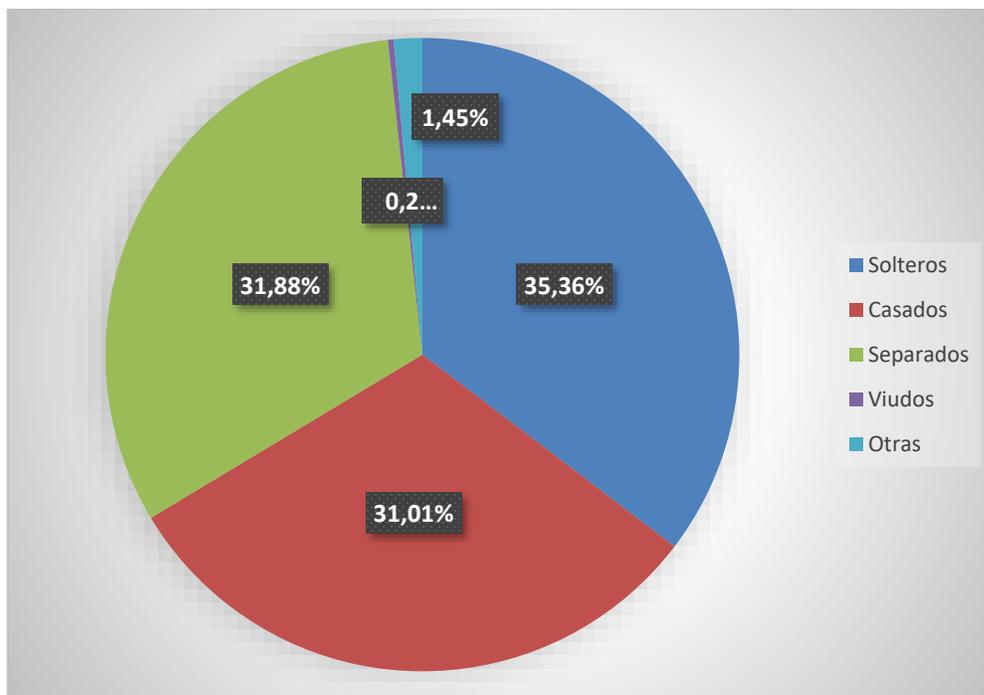
Los diagramas de este tipo se representan sobre un plano delimitado por dos vectores llamados ejes cartesianos. Uno de esos vectores es conocido como eje x o abscisa, y es la horizontal del plano; el eje y , conocido como ordenada, es el eje vertical del plano. En todo gráfico de esta naturaleza, los ejes cartesianos deben estar debidamente señalados. En este caso se observa que las categorías de las variables (solteros, casados, etc.) se han puesto sobre el eje x . El eje de las y contiene los valores de frecuencia absoluta. De este modo, la altura de la barra en un punto dado, representa la cantidad de casos que se han contado en la categoría.

Este arreglo visual de la tabla de frecuencia, llamado diagrama de barras, es particularmente útil para visualizar la distribución de frecuencia entre las categorías. Dado que se trata de una variable de tipo nominal, las barras están separadas entre sí, denotando que no existe continuidad entre las categorías. Se observa entonces que la

población analizada consta mayoritariamente de individuos solteros, casados y separados.

Diagrama de sectores

El diagrama de sectores o gráfico de tortas es una alternativa al diagrama de barras, cuando se utilizan porcentajes. En este caso, el área total de una circunferencia, o bien los 360 ° de giro sobre la misma, representa el 100% de los casos. Los porcentajes propios de cada una de las categorías de la variable, se equiparan al tamaño de una porción o sector del gráfico. El gráfico correspondiente a la columna porcentaje de la variable estado civil, se presenta de la siguiente forma:

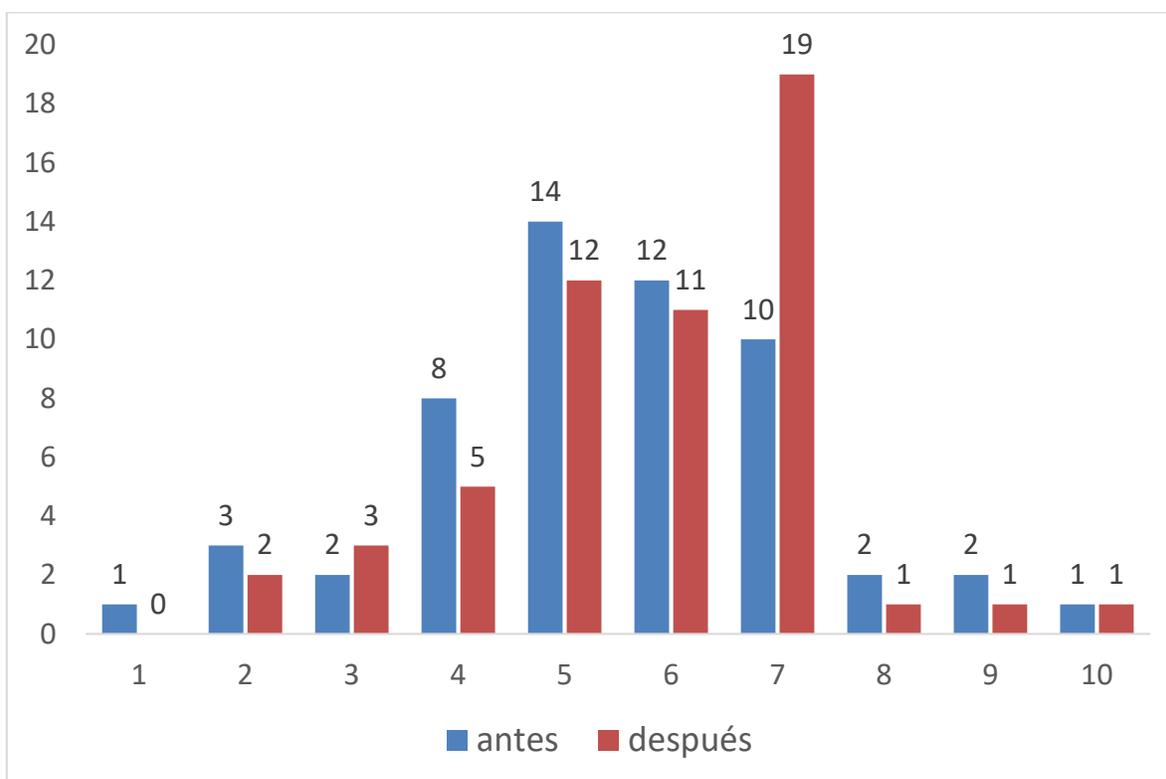


Nótese que el diagrama de sectores prescinde de los ejes cartesianos, dado que la información principal que se usa para el análisis es el tamaño relativo de cada porción de la circunferencia, la cual expresa el porcentaje correspondiente a la categoría de la variable.

Diagrama de barras agrupadas

Este tipo de gráfico tiene las mismas propiedades que el gráfico de barras, pero en él es posible agrupar barras que representen distintas mediciones de una misma

variable. El siguiente gráfico ejemplifica los datos de la tabla donde se midió mediante una escala ordinal, el nivel de cohesión del grupo, antes y después de haber aplicado una serie de actividades que supuestamente la favorecían. Nótese que el eje vertical o eje y, contiene la frecuencia absoluta, pero ahora el eje horizontal o eje de las x, contiene las categorías de la variable ordenada. La medida en que las actividades propuestas han contribuido a aumentar la cohesión grupal, se debe interpretar a partir de las diferencias encontradas en la altura de las barras, en las dos condiciones en que se realizó la intervención: antes – después. Conviene señalar que cada vez que se realizan comparaciones, en mediciones nominales u ordinales, se utiliza preferentemente este tipo de gráficos.



Según se esbozó en la interpretación dada para la tabla de frecuencias, se aprecia que las actividades tendientes a favorecer la cohesión grupal, produjeron como resultado una acumulación de casos en la categoría 7 de la escala, sin modificar sustantivamente la tendencia general de la distribución de frecuencias.

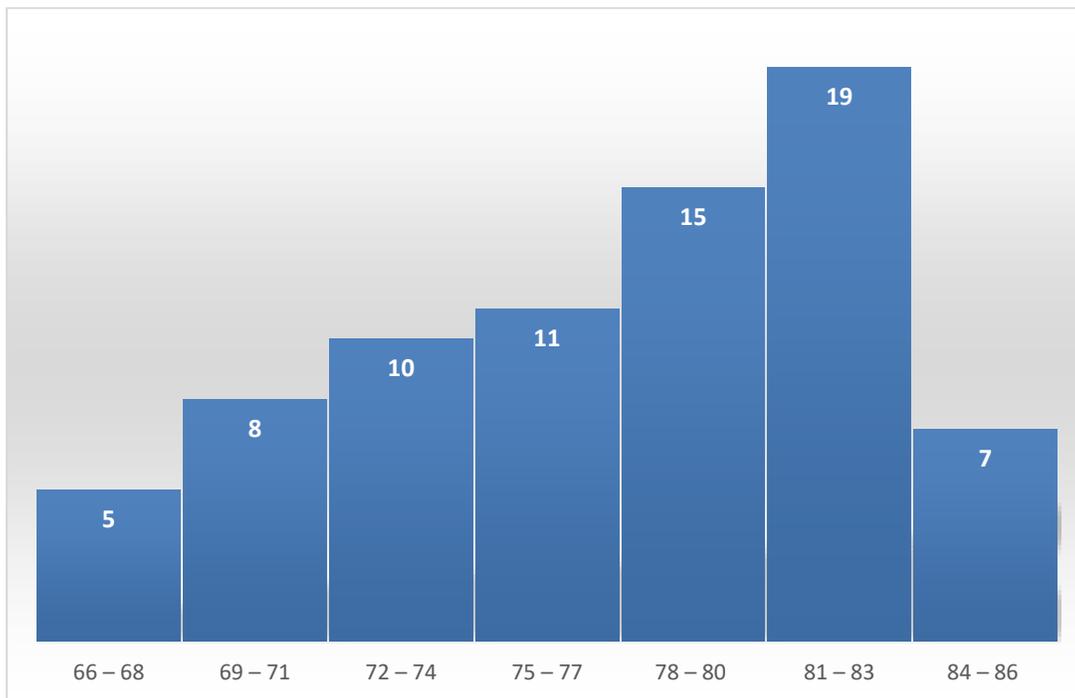
Histograma

Los datos que se muestran en la tabla que corresponden a una variable de tipo métrica, se representan con gráficos similares a los anteriores (diagrama de barras),

pero en los que las columnas son adyacentes entre sí. Ello evidencia que los intervalos de la variable son categorías de un conjunto de datos continuo o discreto. Generalmente, una primera indicación del tipo de datos que se está presentando, viene dado por el tipo de gráfico que se utiliza, el histograma y otros gráficos similares indican claramente el uso de variables métricas. El ancho de las barras de un histograma y su cantidad dependen de la amplitud del intervalo utilizado. La cantidad de intervalos usados para la representación gráfica deben ser suficientes para que el histograma refleje la tendencia general de la distribución de frecuencias. En este sentido, si se usan intervalos muy amplios el recorrido de los valores de la variable queda contenido en muy pocos intervalos, y el histograma resultante tendrá pocas barras, lo cual puede dificultar la observación de la forma de la distribución de frecuencias. Lo mismo ocurre si el intervalo es estrecho y quedan varios de ellos con frecuencia igual a cero.

El histograma para los puntajes de la prueba de madurez lectora, utiliza un intervalo estrecho, pero suficiente para graficar la tendencia general de los datos dado que contiene solo un intervalo con frecuencia cero al comienzo de la distribución.

Prueba de madurez visoespacial: Histograma



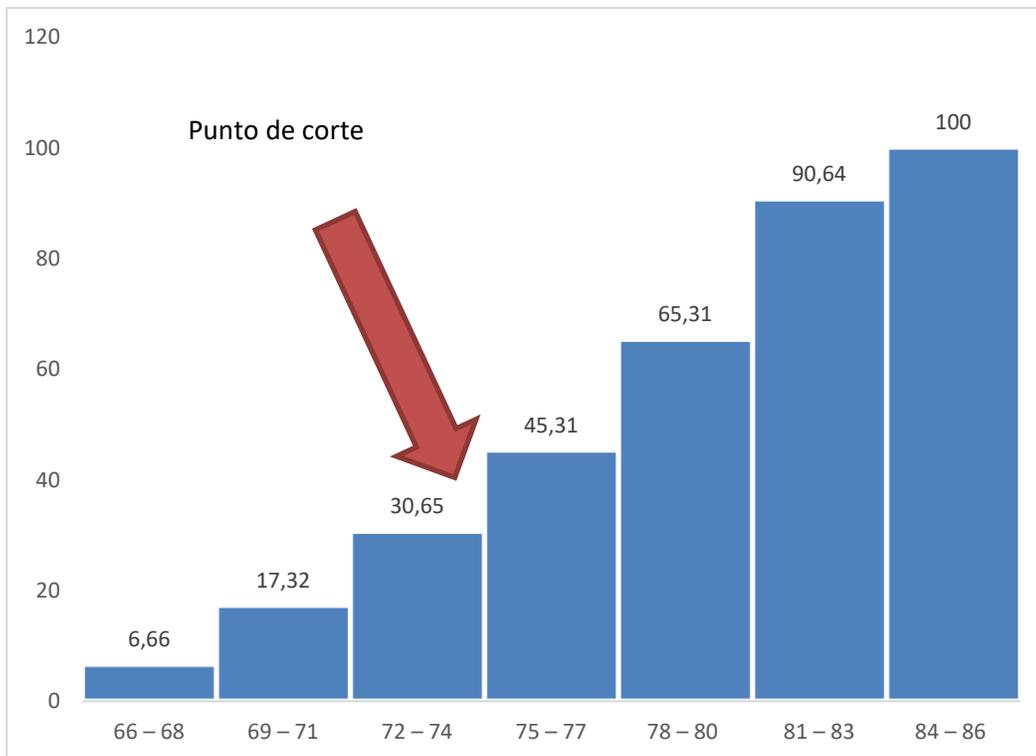
Aquí también se tiene que el eje de las y (vertical), representa la frecuencia absoluta, que en este caso se consignan en cada barra del gráfico; el eje de las x (horizontal),

representa los intervalos de los valores en los que se agruparon las mediciones originales. Por lo tanto, al igual que en el diagrama de barras, la altura de las mismas corresponden a la frecuencia observada en ese intervalo.

Dado que esta es una variable de tipo métrica, es posible realizar un histograma con los valores de frecuencia absoluta o porcentual acumulados. En ambos casos, el histograma representa la progresiva acumulación de casos a través de las categorías de la variable, y ello tiene importancia cuando es posible establecer algún punto de corte teórico para la variable que se está analizando. Recuérdese que en el ejemplo de aplicación de la prueba de madurez lectora, se trazaron límites en los puntajes para favorecer una interpretación cualitativa de la misma, que se reproducen en el siguiente cuadro:

Puntajes de prueba	Interpretación
60 o menos	El niño no posee aún las aptitudes suficiente para el aprendizaje del alfabeto
74 o menos	El niño posee aptitudes limitadas para el aprendizaje del alfabeto
75 o más	El niño posee aptitudes suficientes para el aprendizaje del alfabeto

Estos valores teóricos pueden insertarse en el gráfico para subrayar su importancia al momento de la interpretación.



Siendo un puntaje de 74 el valor por encima del cual se considera que el niño ha desarrollado adecuadamente las habilidades visoespaciales para la lectura, el histograma de porcentajes acumulados muestra que por debajo de ese valor se encuentran el algo más del 30 % de la muestra (23 niños).

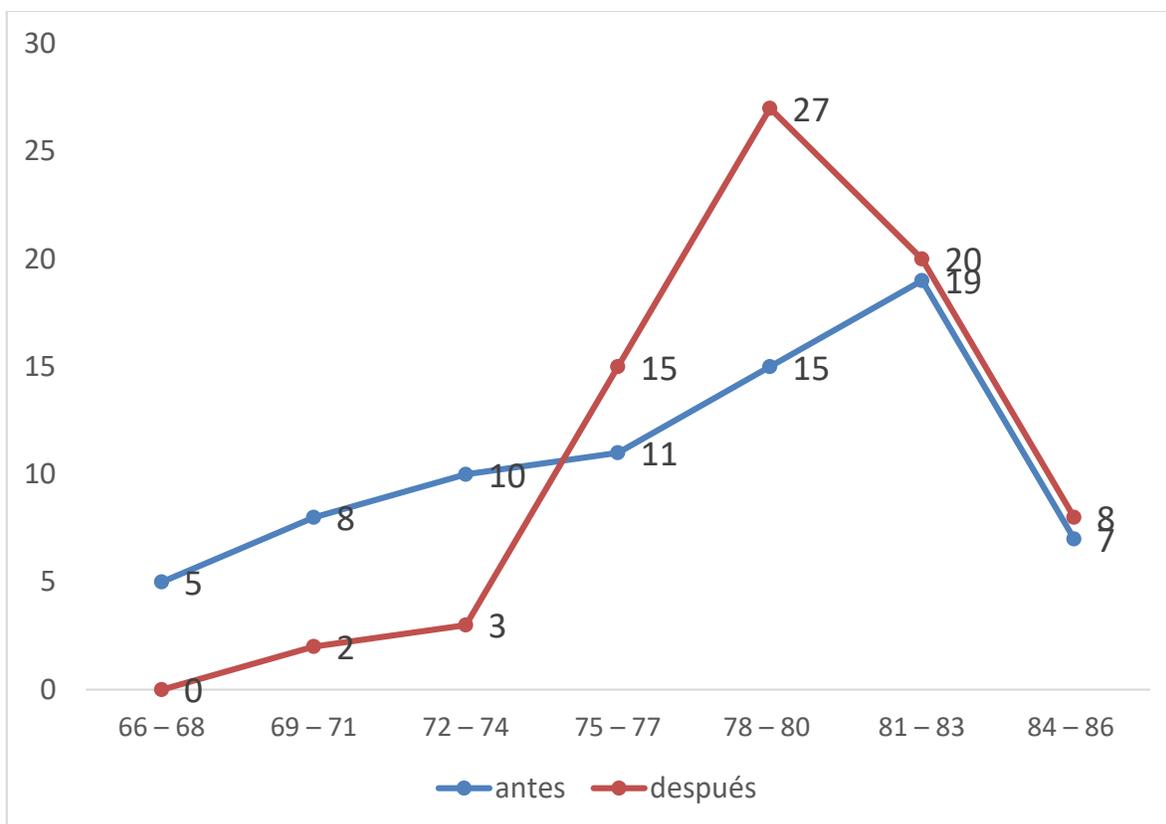
Polígonos de frecuencia

Anteriormente vimos que podía usarse el diagrama de barras agrupado para hacer comparaciones. Con los histogramas una comparación de ese tipo se hace difícil cuando se tiene varios grupos. Sin embargo, la altura de la barra puede reemplazarse por un punto (recuérdese que la altura representa el valor de frecuencia de la categoría), y luego los puntos obtenidos pueden unirse entre sí para generar un polígono. Al reemplazar el histograma por el polígono, es posible comparar la distribución de frecuencia de dos o más grupos en una misma variable métrica.

Para este caso, ampliaremos el ejemplo anterior sobre la prueba de madurez lectora. Supongamos que la muestra de escolares a las que se les aplicó la prueba de madurez visoespacial, realizaron actividades tendientes a estimular el reconocimiento de las letras del alfabeto. Los maestros deciden reevaluar a los niños tres meses después y verificar si hubieron cambios en los puntajes de prueba. Explícitamente, se espera que aquellos niños con menor puntajes alcancen una puntuación mayor, pero

no se anticipan mayores cambios en aquellos escolares que ya poseían la madurez visoespacial suficiente. Un resultado de este hipotético experimento se muestra en el siguiente polígono de frecuencia. Nos vamos a centrar en el puntaje 74 que es el punto de corte que separa los escolares que poseen aptitud suficiente de aquellos que tienen aptitudes limitadas.

Polígono de frecuencia: prueba de madurez visoespacial para la lectura.

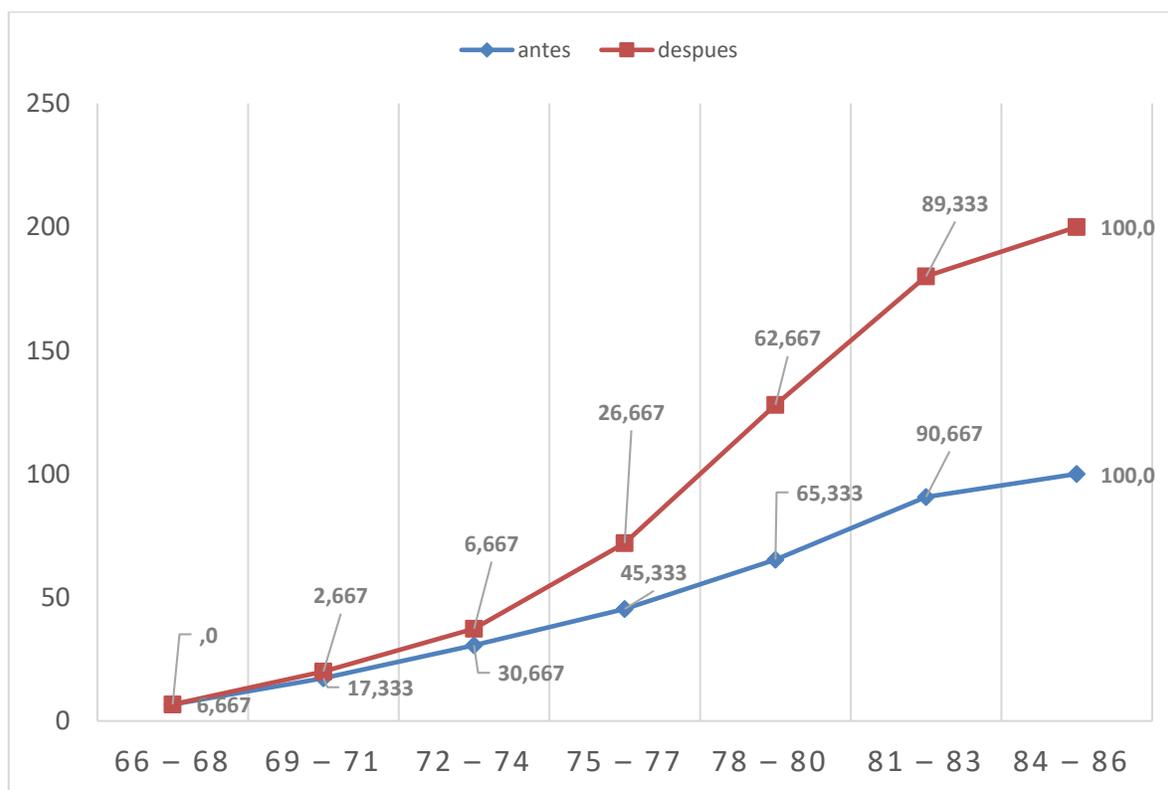


Lo primero que apreciamos en estos datos es que tras la intervención de los docentes, no se verifican casos con puntajes en el rango 66 - 68. La cantidad de casos con puntajes iguales o menores a 74 es ahora de 5 escolares (es la cantidad de alumnos en los intervalos anteriores a 75-77), cuando antes era de 23. Con esos sencillos datos es posible afirmar que la intervención realizada ha resultado beneficiosa para aquellos escolares con aptitudes visoespaciales menos desarrolladas al comienzo del aprendizaje de la lectura. Como era de esperar, aquellos alumnos que alcanzaron los mayores puntajes de prueba antes de la intervención, no modificaron sustantivamente su posición en la prueba tres meses después. Esto era de esperar debido a que las actividades desarrolladas no se proponían estimular a aquellos con las aptitudes

visoespaciales suficientes, sino a aquellos con menores niveles de esta habilidad.

Gráfico de ojiva

El gráfico que se presenta a continuación corresponde a los valores observados de la variable puntajes obtenidos en una prueba de lectura. En él, cada punto representa el valor observado de la variable y la frecuencia que se acumula con cada valor. Este tipo de gráfico es utilizado para representar el total acumulado a través de los valores ascendentes de la variable.



En esta gráfica vemos el porcentaje acumulado de casos según los puntajes en la prueba de madurez visoespacial. Se puede ver claramente que el gráfico de ojiva es una variante del polígono de frecuencia. En la situación "antes", cuando no se ha realizado ninguna intervención, vemos que la ojiva tiene una pendiente menos pronunciada, comparada con la situación "después", momento en el que se realizó la intervención con actividades para estimular el aprestamiento visoespacial. La diferencia entre ambas curvas u ojivas se debe a la celeridad con la que se acumulan los casos. Cuantos más casos se acumulan en cada intervalo, más pronunciada la

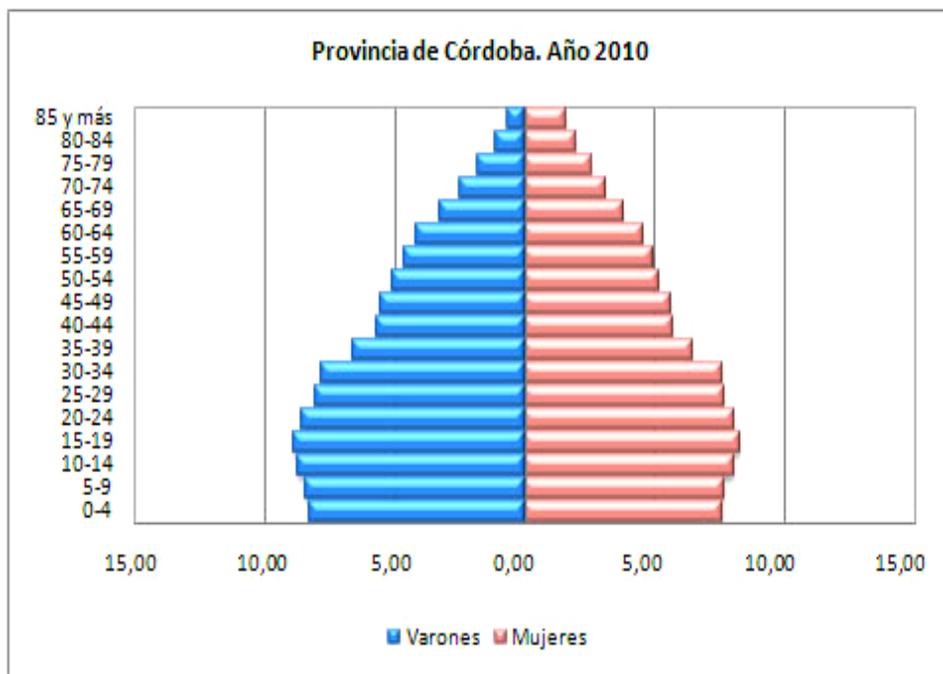
pendiente, la distancia entre ambas curvas marca la diferencia entre lo ocurrido antes de la intervención y posterior a la misma.

Otros tipos de gráficas

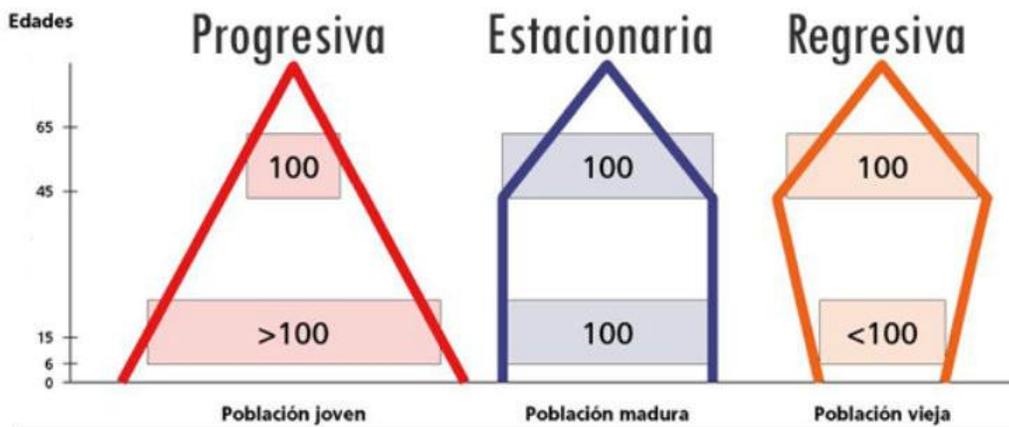
Los gráficos mostrados hasta aquí se corresponden con el tipo de variables que hemos analizado, pero en estadística existen muchos tipos de gráficas diferentes dependiendo del tipo de información que se quiera transmitir, del énfasis que se quiera dar a ciertos aspectos de los datos, si se utilizará con fines científicos o de divulgación. Muchas publicaciones especializadas requieren que los autores utilicen ciertos estilos y formatos en la presentación de datos. Por lo dicho, en los apartados siguientes mostraremos algunos gráficos frecuentemente utilizados en reportes estadísticos.

Pirámides de población

La pirámide de población es una gráfica muy utilizada en epidemiología y en estudios demográficos. Se trata de la superposición de dos histogramas, uno que corresponde a varones y el otro a mujeres. En el eje de ordenadas se representan por intervalos las edades de los individuos, y en el eje de la abscisa el porcentaje de población en ese rango de edad. Dado que la cantidad de población disminuye con la edad, lo común es ver este gráfico en forma de pirámide. Cuando esto sucede, el exponente de crecimiento poblacional es positivo, dado que hay más nacimientos que defunciones, y la población joven sobrepasa a la añosa. Pero suele suceder que en algunas regiones o países, esa pirámide se invierte y el pico es mayor que la base; en tal caso el exponente de crecimiento poblacional es negativo. El gráfico que sigue corresponde a la población de la Provincia de Córdoba; en él se aprecia que el exponente de crecimiento poblacional es positivo por tanto la pirámide es más ancha en su base. En el eje de las x la referencia está dada por la cantidad de varones y mujeres en miles de personas. En el eje de las y se hace referencia a la franja de edad correspondiente. Bajo condiciones de crecimiento normal vegetativo de la población, la distribución de la población de varones y mujeres tiende a ser simétrica. Un dato relevante es la pérdida de simetría del gráfico a partir de la franja etaria de 65 a 69 años, que mostraría que las mujeres tienden a ser más longevas que los varones.

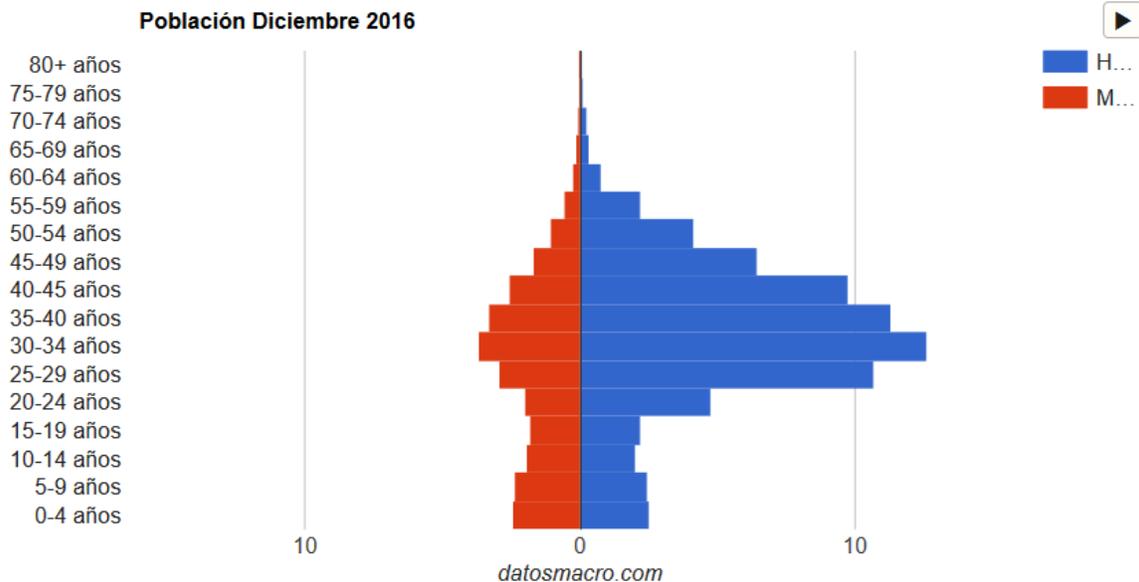


Los estudios de poblaciones diferencian al menos tres tipos de pirámides de población. Un tipo de pirámide llamada progresiva, tiene base ancha y se va reduciendo paulatinamente hacia la cúspide. Es consecuencia de una alta tasa de natalidad y con población mayormente joven. Muchos países latinoamericanos presentan este tipo de crecimiento poblacional. También existen las pirámides estacionarias donde se manifiesta que la tasa de nacimientos es menor y está equilibrada con la población joven, esto da cuenta que la natalidad y mortalidad se han mantenido sin variaciones significativas durante un periodo de tiempo largo. Esta pirámide se considera el paso intermedio entre la pirámide progresiva y la regresiva. La pirámide regresiva es más ancha en los grupos superiores que en la base, debido al descenso en la natalidad y al envejecimiento de su población. A medida que descienden los nacimientos y la población se hace añosa la pirámide toma un aspecto invertido. Actualmente, en varios países europeos se verifica este tipo de crecimiento poblacional. En el siguiente gráfico podemos apreciar las diferencias mencionadas.



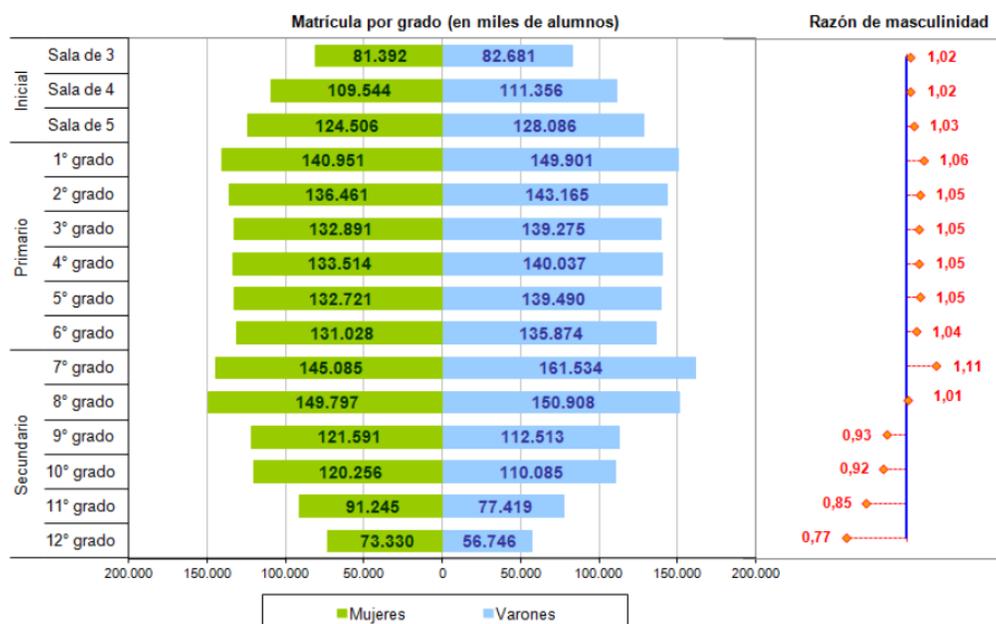
Cabe destacar que muchos fenómenos sociales pueden intervenir para que la forma esperada de la distribución de la población pierda su armonía entre varones y mujeres. Durante el siglo XX las dos guerras mundiales alteraron la relación entre varones y mujeres en las pirámides de población europeas. Actualmente, son las migraciones en masa las que están produciendo efectos de distorsión en el crecimiento poblacional. En el gráfico que se muestra más abajo se aprecia una cantidad mucho mayor de hombres jóvenes con relación a mujeres en los Emiratos Árabes Unidos, especialmente en su capital Abu Dabi. El gráfico que se muestra a continuación es del total del país y la causa del crecimiento de población masculina se debe a la migración para los trabajos de construcción en la capital del país.

Pirámide de población para Emiratos Árabes Unidos.



Las pirámides de población también pueden utilizarse en educación para mostrar la cantidad de estudiantes por nivel según el sexo. En la siguiente gráfica se muestra la tasa neta de escolarización de la provincia de Buenos Aires para el año 2010. Un dato que suele incluirse en estas gráficas es la razón de masculinidad, la cual resulta de dividir la cantidad de varones por mujeres. Si se hubiera realizado la operación inversa se tendría la razón de femineidad. Lo que se observa en este gráfico es que, para el año reportado, la razón de masculinidad se invierte en el tramo final del trayecto educativo. La razón es simplemente que para esa franja etaria hay más mujeres que varones.

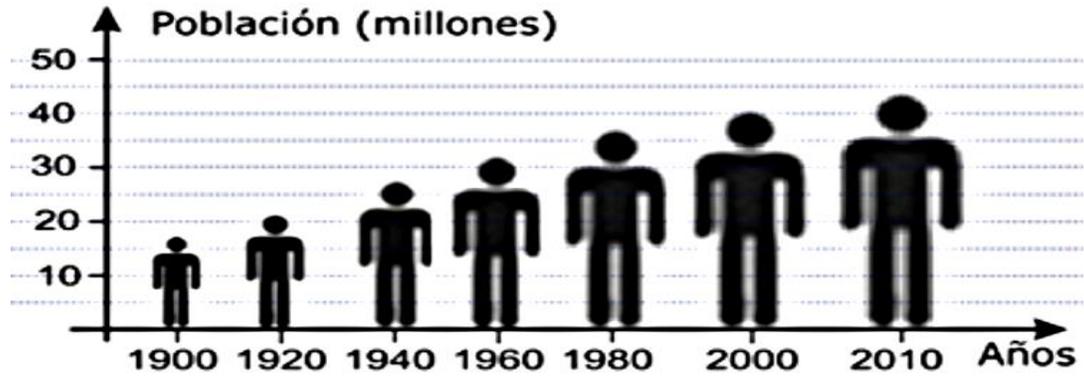
Matrícula por edad simple y condición de edad, por sexo. Educación común. Año 2010. En miles de alumnos



Fuente: Procesamientos propios sobre datos de DiNIECE-ME. Relevamiento Anual de Matrícula y Cargos. Anuario estadístico 2010

Pictogramas

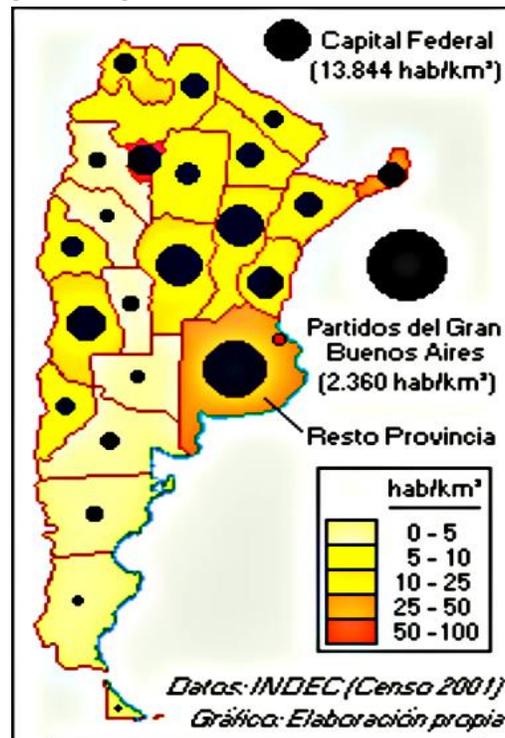
El pictograma no es otra cosa que un diagrama de barras, que incorpora figuras que simbolizan a la variable en estudio. No suele presentar mayor información que el diagrama de barras, pero en ocasiones es más sencillo de interpretar por su recurso iconográfico. Nótese cómo en el gráfico que presentamos a continuación, aparece claramente expresado el crecimiento de la población argentina mediante el uso de una figura humana en reemplazo de una barra. Una cuestión que siempre debe tenerse en cuenta a la hora de interpretar este tipo de gráfica, es la cantidad de individuos (o unidades de análisis) que se representan en cada pictograma. En este caso cada figura representa a millones de personas.



En el año 2010 en argentina se contabilizaron 40.091.359 personas

Mapas

Como se habrá notado en el pictograma, la representación gráfica es aproximativa y sus valores no están rigurosamente escalados. Cuando se presenta esta dificultad y se requiere precisar información, es posible utilizar varios recursos simultáneos, tal como se muestra en el siguiente gráfico.



Este gráfico utiliza tres recursos: a) el mapa del territorio argentino, b) puntos cuyo tamaño representa la cantidad de población concentrada en cada provincia y c) una

paleta de colores que representa la densidad poblacional. Combinando los tres recursos es posible ver que la provincia más densamente poblada es Buenos Aires, pero la que concentra mayor cantidad de habitantes por Km² es la provincia de Tucumán.

Gráficas Polares

La gráfica polar fue una ingeniosa idea para agregar información al diagrama de sectores, ideada por la enfermera Florence Nightingale durante la guerra de Crimea (1853-1856). En su trabajo en los hospitales de campaña, Florence Nightingale recopiló datos que demostraban que las heridas de guerra eran una fracción menor de muerte de los combatientes, en comparación con enfermedades como el tifus, el cólera y la disentería. Estas últimas eran las principales causas de muerte, debidas principalmente a las deficientes condiciones de higiene hospitalaria. El desafío que enfrentó Florence Nightingale, fue mostrar de manera convincente a las autoridades, la necesidad de cambiar de manera drástica y urgente las condiciones de prácticas quirúrgicas y médicas. El diagrama que se muestra a continuación fue el que cambió el curso de las atenciones recibidas por los soldados durante la guerra.

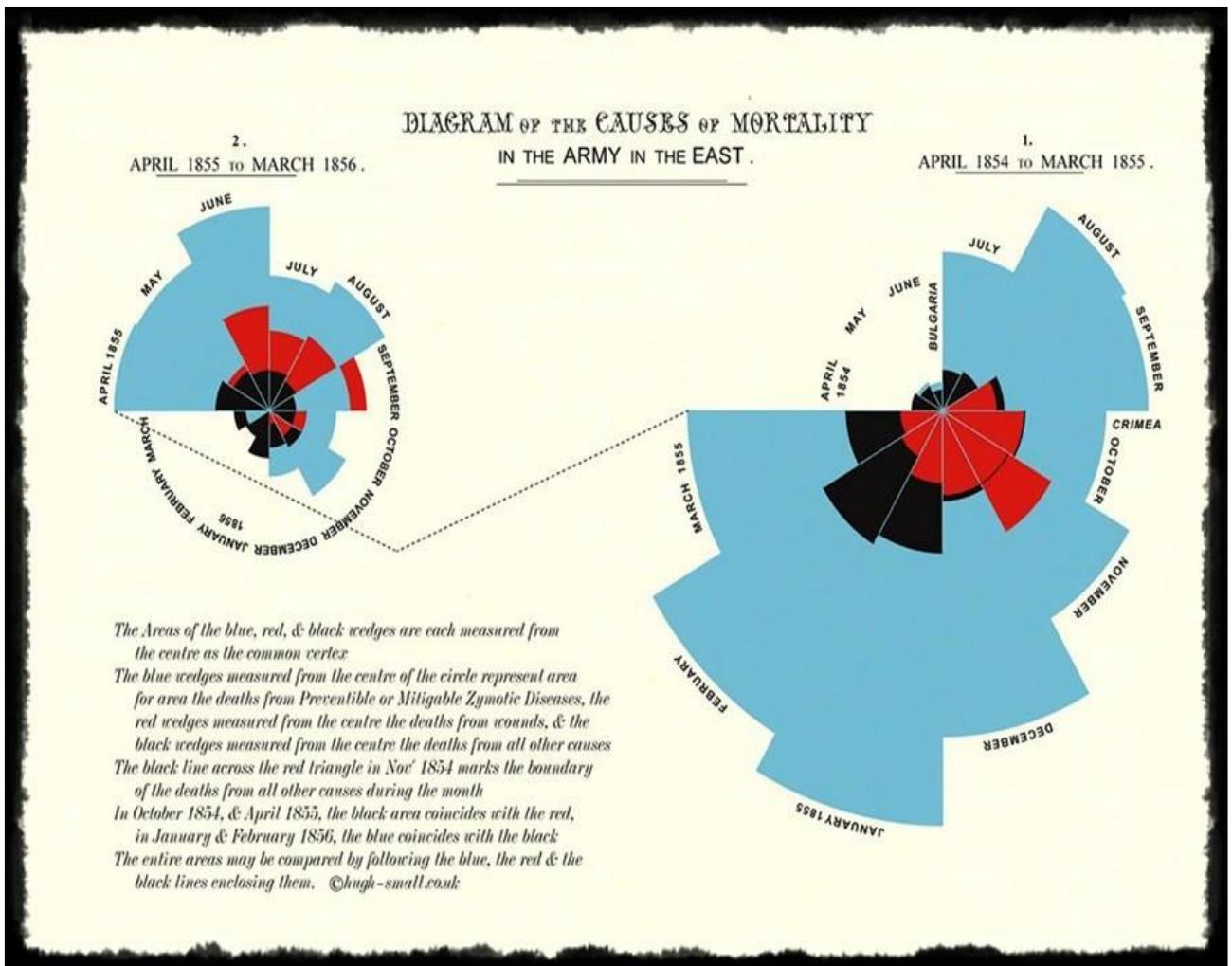
Interesa destacar en este caso que cada sector del gráfico se corresponde a un mes. Luego, los colores indican las causas de muerte de los soldados:

Azul: Muertes por enfermedades infecciosas, desde prevenibles hasta mitigables

Rojo: Muertes por heridas de guerra

Negras: Muertes por otras causas

Nótese cómo el sector azul crece mucho más que el sector rojo alcanzando su pico en enero de 1855; en julio de 1854 ya es posible ver que las enfermedades aventajan en causa de muerte a las heridas en combate. La información presentada por Nightingale fue decisiva para promover el cambio en las prácticas hospitalarias, y su enfoque estadístico convenció a las autoridades militares, al parlamento y a la reina Victoria, para llevar a cabo la reforma. Gracias a ella se mejoraron las condiciones de sanidad y se consiguió reducir la proporción de muerte de sus pacientes. En febrero de 1855 la tasa de mortalidad había descendido del 60% al 42,7% y en primavera ya era del 2,2%.

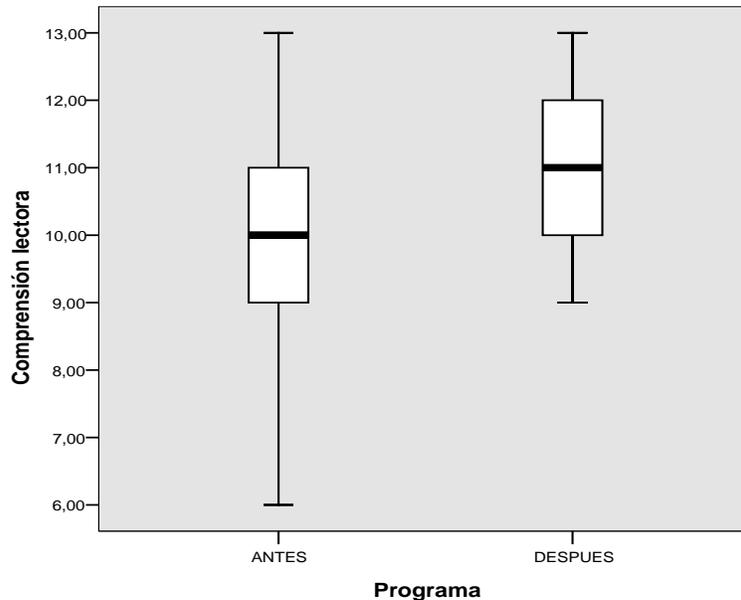


Box Plot o Diagrama de Cajas

El box-plot o diagrama de cajas es una sencilla representación de la dispersión de datos para una variable que permite comparar fácilmente el comportamiento de dicha variable en uno o más grupos. El diagrama exhibe como puntos de referencia un rectángulo (caja), intersectado por una línea. Generalmente, estos puntos representan la posición de los cuartiles de la distribución de datos (eventualmente, la línea media en la caja puede representar el promedio del grupo). De cada extremo de la caja se extienden líneas perpendiculares a la misma y representan en cada caso, la distribución de datos hacia el máximo y el mínimo valor registrado en la distribución de valores de la variable. La interpretación de un diagrama de cajas depende de la posición de los puntos de referencia del diagrama respecto de la distribución de datos observada, es por ello que el diagrama siempre se representa en el marco de la intersección de ejes cartesianos.

Diagrama de cajas y medidas de posición: cuartiles

Para comprender mejor la utilidad de este tipo de diagramas, proponemos el siguiente ejemplo: Una investigadora desarrollo un programa de lectura compartida para aplicar en el aula. Según ella el programa estimula la lectura y mejora la calidad de la comprensión. Para verificar la eficacia del programa, seleccionó 196 escolares, en quienes evaluó la comprensión lectora antes de la aplicación del programa, y después de dos meses de aplicación del mismo. La comprensión lectora se evaluó mediante una prueba estandarizada, con valores que van de 0 a 15, siendo este último el máximo puntaje de la prueba. En el siguiente gráfico se ofrece un diagrama de caja con los resultados obtenidos.



Para comprender los diagramas de cajas primeramente se debe atender a la referencia que proporcionan los ejes cartesianos. En este caso el eje horizontal indica la medición de la variable en dos situaciones diferentes: antes de aplicar el programa y después de aplicado el mismo. El eje vertical es la referencia para la distribución de valores en la variable comprensión lectora. Así, es posible observar que se ha producido un cambio en la variable una vez aplicado el programa, dado que la posición de ambas cajas es diferente.

Una primera interpretación comienza con la posición que ocupa el valor

máximo y mínimo en cada distribución; de este modo se observa una diferencia en el valor mínimo de cada distribución. Tal diferencia es de tres puntos en la variable, e indica un incremento del valor mínimo luego de la aplicación del programa. Por otro lado, se observa que el valor máximo no se ha modificado.

La longitud de la línea inferior que alcanza el rectángulo recibe el nombre de bigote (al igual que la línea superior). Indican la posición en los valores de la variable para $\frac{1}{4}$ de los casos. En otras palabras, la longitud del bigote indica qué valores de la variable han obtenido aquellos casos que se encuentra entre el mínimo valor registrado y el primer cuartil. Como se aprecia la distribución es diferente antes y después de aplicar el programa; antes de aplicarlo $\frac{1}{4}$ de los casos obtuvo valores de 6 a 9 en comprensión lectora, luego de aplicado el programa $\frac{1}{4}$ de la distribución se concentra en los valores 9 y 10. Por tanto se registra un incremento en la comprensión lectora de aquellos individuos que antes habían mostrado el menor rendimiento.

La parte inferior del rectángulo y la línea media que lo divide indica la distancia entre el primer cuartil y el segundo cuartil o mediana. La distancia entre ambos indica los valores obtenidos en la variable por $\frac{1}{4}$ de la distribución. Nótese que la cantidad de casos acumulados entre el mínimo y la mediana corresponde a la mitad de los casos de la muestra. En el diagrama, se observa que la dispersión en ambos grupos, entre el primer y segundo cuartil, es la misma, sin embargo, los valores de la variable son superiores luego de la aplicación del programa. Dicho en otras palabras, antes de la aplicación del programa $\frac{1}{4}$ de la distribución de casos obtuvo valores entre 9 y 10; luego de aplicar el programa $\frac{1}{4}$ de la distribución alcanza valores entre 10 y 11. La situación es similar para todos los casos que se encuentran entre el segundo y tercer cuartil (que corresponde también a $\frac{1}{4}$ de los casos).

Finalmente, y dado que no se ha modificado el máximo valor de la variable, se observa que luego de la aplicación del programa la dispersión es menor. En otras palabras, antes de la aplicación del programa el cuarto de los casos con mayor rendimiento en comprensión lectora alcanzaba puntaje de prueba entre 11 y 13; luego de la aplicación este grupo alcanza valores concentrados en 12 y 13. En conclusión, la comparación mediante el diagrama de cajas de la situación antes y después de la aplicación del programa, muestra que este resultó efectivo en la muestra analizada, para mejorar el nivel de comprensión lectora.

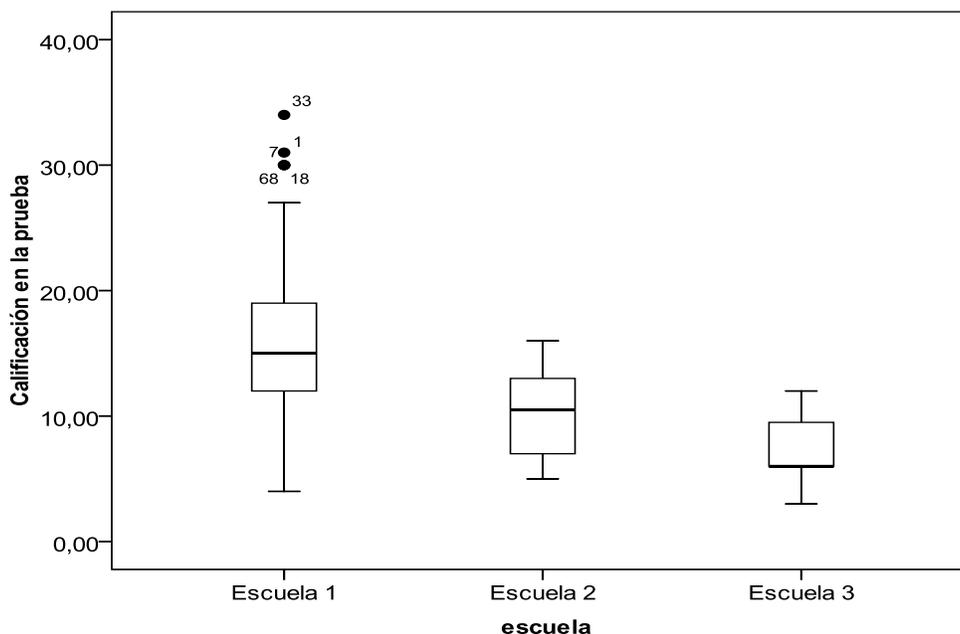
Diagrama de cajas – casos con valores atípicos o extremos

Para determinar que una distribución de datos contiene valores extremos, se toma en consideración una medida proporcional del rango intercuartil, mediante la cual se determina la extensión de los máximos y mínimos (bigotes de la caja). El rango

intercuartilar se puede calcular de diferentes maneras, aquí consideraremos las dos más usadas:

- a) Rango intercuartilar (RI): $Q3-Q1$
- b) Rango semiintercuartilar (RSI): $Q3-Q1/2$

Para desarrollar este tema proponemos el siguiente ejemplo: aplicando una prueba estandarizada para evaluar razonamiento numérico, se compararon los puntajes obtenidos por tres muestras de escolares del segundo ciclo de EGB de tres escuelas diferentes. El siguiente diagrama de cajas muestra los resultados obtenidos en las evaluaciones realizadas.



En este ejemplo se tiene las referencias de cada caja en el eje horizontal (escuela), y la distribución de los valores de la variable en el eje vertical (puntaje en la prueba de razonamiento numérico). En lo esencial la interpretación del diagrama es el mismo que hemos desarrollado anteriormente, con la excepción de que los bigotes no se extienden del máximo al mínimo puntaje, sino que representan una proporción del RSI. Eso hace que en la escuela 1 se observen casos atípicos. En esa muestra existen cinco casos atípicos acumulados en el extremo superior de la distribución de la variable.

Considere la posición de $Q1$ y $Q3$ en la distribución de valores de la variable razonamiento numérico de la escuela 1 en el gráfico, se tiene que $Q1=11$ y $Q3=19$ (los valores son aproximados). Por tanto, se tiene que $RI= 8$; $RSI= 4$. Los casos extremos se identifican cuando muestran valores que exceden la longitud de los bigotes, y esta longitud es una proporción de RI o de RSI . Digamos que si tomamos 1.5 del valor de RSI , el recorrido de los bigotes se extiende desde 19 a 25 en el límite superior, y de 5 a

11 en el límite inferior. Estos valores surgen de los siguientes cálculos:

$$1.5 \times \text{RSI} = 1.5 \times 4 = 6$$

$$Q3 + (1.5 \times \text{RSI}) = 19 + 6 = 25$$

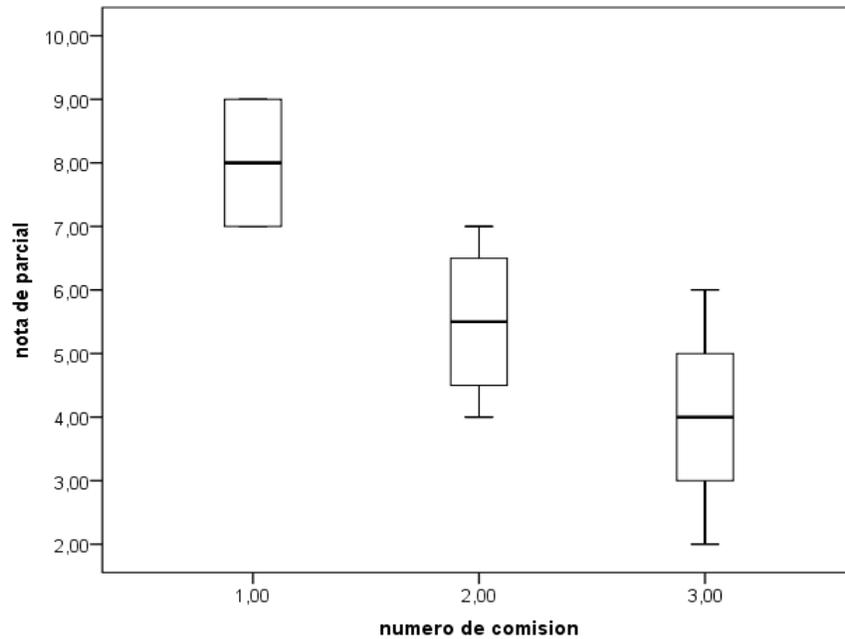
$$Q1 - (1.5 \times \text{RSI}) = 11 - 6 = 5$$

Se aprecia que en una distribución de datos cuyos límites máximos y mínimos son una proyección de un valor teórico (en este ejemplo es 1.5 del valor de RSI), pueden existir valores que exceden ese límite. A tales casos se los denomina atípicos o extremos. En las gráficas estos casos aparecen marcados con puntos y apropiadamente identificados. De este modo, sabemos que los casos 68 y 18 son atípicos puesto que obtuvieron un puntaje de 30; del mismo modo, el caso 7 y 1 obtuvieron un valor de 31 y el caso 33 obtuvo un valor de 35 aproximadamente. Nótese que los casos atípicos aparecen solo en la escuela 1, la cual es la que muestra la mayor dispersión en los puntajes de la prueba de razonamiento numérico.

Otro aspecto interesante presente en el gráfico aparece en la distribución de la escuela 3, en la cual se observa que no hay variación entre Q1 y Q3, lo cual queda reflejado visualmente en que no existe separación entre la parte inferior del rectángulo y la intersección que marca la mediana. En otras palabras, la mediana (Q2) y el primer cuartil (Q1) tienen el mismo valor, indicando que 1/4 de la muestra no evidencia variación en el puntaje de la prueba de razonamiento numérico (todos los casos obtuvieron una calificación de 8).

Diagrama de cajas – restricción de la variabilidad

Así como es posible identificar casos extremos en diagrama de cajas, también se puede identificar la situación en que los casos se concentran en un valor determinado de la variable. Anteriormente vimos que la restricción de la variabilidad hace que coincida el valor de mediana con el primer cuartil. En el gráfico que sigue, la restricción de la variabilidad se expresa en la concurrencia del valor máximo y mínimo con el tercer y el primer cuartil respectivamente. Veamos el siguiente ejemplo: en la materia Sociología se dividieron a los alumnos en tres comisiones diferentes. Los docentes supusieron que el orden de las comisiones del parcial no alteraría el rendimiento de los estudiantes. Atendiendo a los resultados de los exámenes se elaboró el siguiente diagrama de cajas.

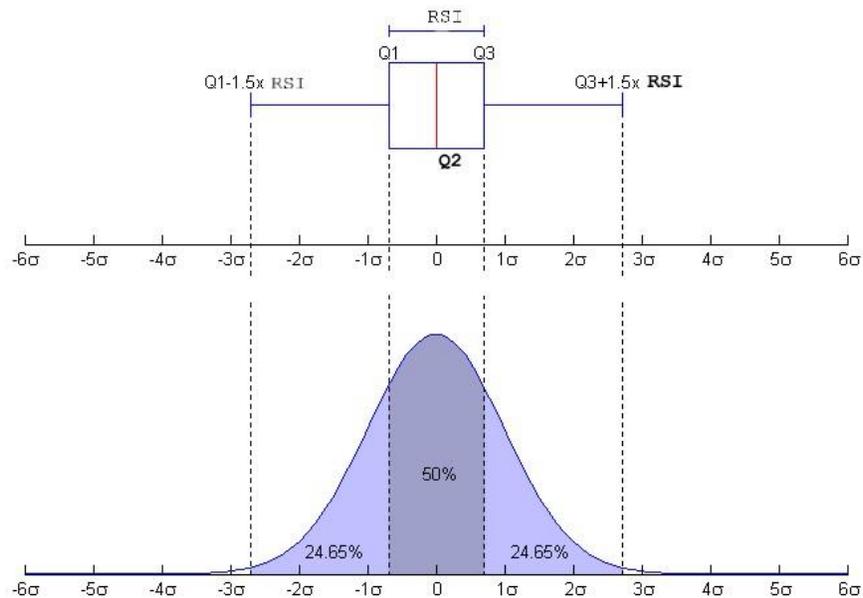


Tal como se observa, el rendimiento de los alumnos muestra un patrón que los diferencia claramente. Sobresalen dos aspectos, el primero es que el rendimiento es menor a medida que se pasa de la comisión 1 a la tres, el segundo es que, siguiendo ese mismo orden, los resultados de los exámenes muestran mayor variabilidad. En la comisión 1 se observa una particularidad que es la restricción en la variabilidad. En esta comisión los alumnos obtuvieron 7, 8 y 9. Más precisamente, el diagrama muestra que un tercio del total de alumnos obtuvo un siete, otro tercio obtuvo un ocho y otro más un nueve, completándose así el total de la muestra de alumnos de esa comisión. En una situación como la descrita, el máximo y el mínimo de la distribución no se diferencian del primer y el tercer cuartil. En otras palabras, no existe variabilidad más allá de los cuartiles.

Diagrama de cajas en distribuciones sesgadas

Más adelante veremos una distribución muy utilizada en ciencias sociales que permite describir el comportamiento de muchas variables, se trata de la distribución normal. Una de las principales características de la distribución normal teórica es su simetría, y esta propiedad es importante porque es posible determinar (numéricamente o gráficamente) cuándo existen sesgos en una distribución empírica de datos. Los sesgos son indicadores de que los valores de la variable tienden a acumularse hacia alguno de los extremos de todos los valores posibles de la variable en cuestión. No abundaremos en más detalles pues profundizaremos este tema en el capítulo dedicado a distribución normal. En este apartado destacamos que existe una relación entre la forma de una distribución normal teórica y el diagrama de cajas. La siguiente figura

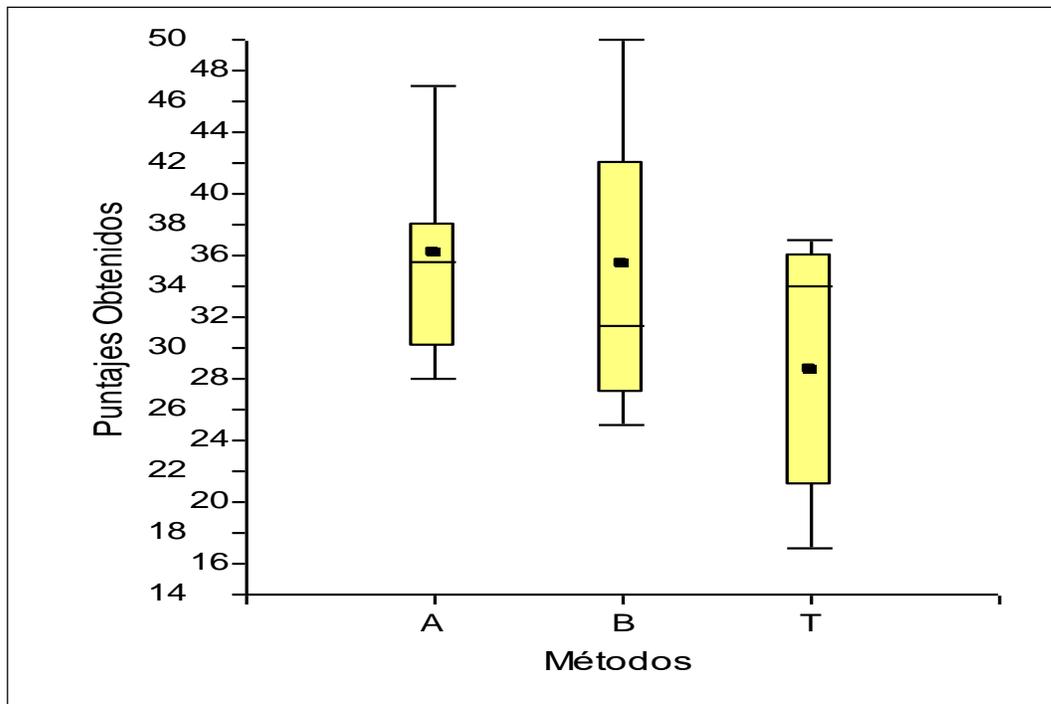
muestra esa relación.



El diagrama de cajas que le corresponde a una distribución normal teórica, muestra que las extensiones de los bigotes de la caja son proporcionales a las extensiones de la posición de los cuartiles. Dicho en otras palabras, la distancia entre el valor mínimo y $Q1$ es la misma que entre el valor máximo y $Q3$; además la distancia entre $Q2$ y $Q1$, es la misma que la que existe entre $Q2$ y $Q3$. La distribución descrita es “ideal”, y sirve de punto de referencia teórico para comparar distribuciones empíricas.

Por lo dicho, para saber si una distribución de datos empíricos muestra algún sesgo, solo basta agregar a un diagrama de cajas una medida complementaria de la mediana, que es la media aritmética y verificar además las distancias entre las distintas medidas de posición. Entonces, si la media aritmética y la mediana se apartan entre sí lo mismo que las medidas de posición, nos informa de la tendencia del sesgo de una distribución. Una aclaración importante en este punto es que no debe entenderse que las distribuciones segadas son “malas distribuciones” puesto que la forma del sesgo informa cómo están ordenados los datos de una variable empírica.

Para completar lo que acabamos de mencionar, veamos un ejemplo: en una escuela se aplicaron tres métodos diferentes para la enseñanza de las matemáticas. Luego de dos meses, los niños que aprendieron con cada método fueron evaluados con una prueba estandarizada cuyo puntaje máximo alcanza 50 puntos. Los resultados de los grupos en la evaluación se muestran en el siguiente gráfico de cajas. En esta gráfica se aprecia que la mediana está representada por una línea transversal a la caja y la media aritmética por un punto en el interior de la caja.



Para este ejemplo utilizaremos las siguientes referencias: el método A y B son los métodos nuevos, los cuales son comparados con el método tradicional (T: Método Tradicional). No analizaremos en detalle la complejidad de los gráficos solo su asimetría. El método A muestra la particularidad de que Q2 y Q3 están más próximos entre sí, que Q2 y Q1, por tanto, un tercio de los casos se ha concentrado en los valores 36, 37 y 38. Por otro lado, los bigotes son diferentes, el inferior es considerablemente más pequeño que el superior, por lo tanto, hay mayor dispersión de casos en los valores altos de la variable. El método B ofrece una distribución diferente y se observa que el sesgo está marcado por los valores que están por encima de la mediana. Esto lo sabemos porque la media aritmética está por encima de la mediana, la distancia entre Q2 y Q3 es mayor que entre Q2 y Q1 y porque el bigote superior es más largo que el inferior. Considerando todos estos aspectos del gráfico concluimos que los que participaron en este método tienden a dispersarse más hacia los valores altos de la variable. Por último, el método tradicional es una imagen casi especular del método B, solo que la tendencia es hacia los valores bajos de la variable. En conclusión, es posible afirmar que los métodos que mejores resultados ofrecen son el A y el B en comparación con el método tradicional. Comparando los métodos A y B respectivamente, el primero muestra una ventaja relativa respecto del segundo.

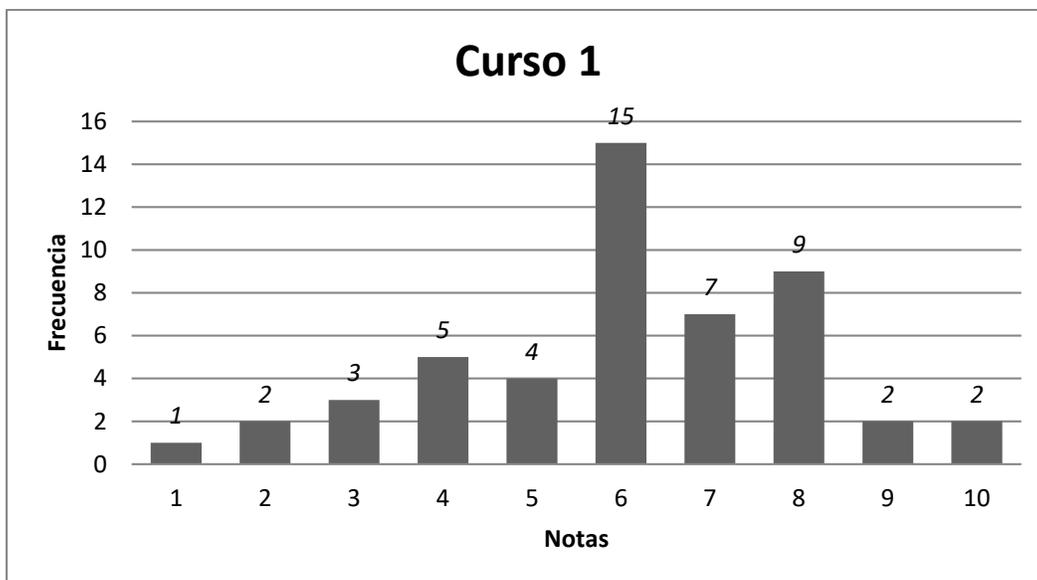
Capítulo 4

Medidas de tendencia central

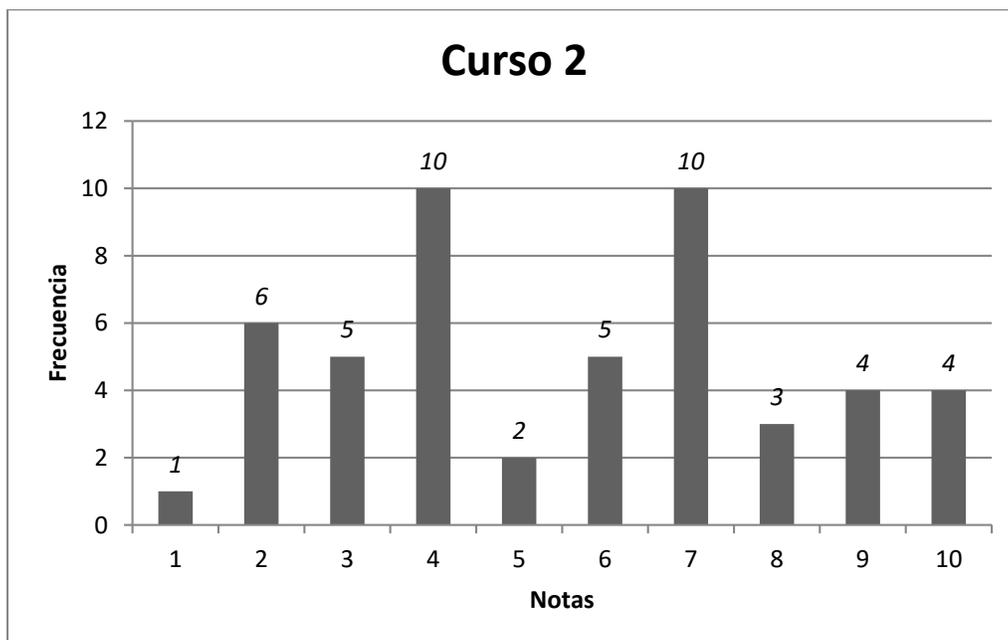
Si se dispone de un conjunto de datos, es posible utilizar algunas medidas que permitan resumir en unos cuantos números significativos, las principales características del conjunto completo. Estos números, cuando son obtenidos de una muestra reciben el nombre de estadísticos y cuando se calculan en una población se denominan parámetros. Algunos estadísticos muestran la posición central del conjunto de datos, por lo cual se reciben el nombre de medidas de tendencia central. Repasaremos aquí las medidas más comúnmente utilizadas.

Modo

También recibe el nombre de moda, y se define como la categoría de la variable que concentra la mayor frecuencia. Si en un conjunto de datos, existe un valor que es el más común, esto es, que se repite más veces y por tanto es el más frecuente, se dice que la distribución de esos datos es unimodal. La mayoría de las distribuciones de datos son unimodales, pero en ocasiones puede que existan dos valores que se han repetido igual número de veces, siendo estos los más frecuentes. En tal caso, la distribución tiene dos modos, y por tanto es bimodal. En el caso de que existan varios valores que concentran las mayores frecuencias de la distribución de datos, se tiene una distribución multimodal. El modo es la medida más sencilla de caracterizar una distribución de datos. Ejemplo: un docente toma un examen de matemáticas en tres cursos de 50 alumnos. Luego realiza una gráfica con la distribución de los resultados, identificando en cada caso el valor modal.

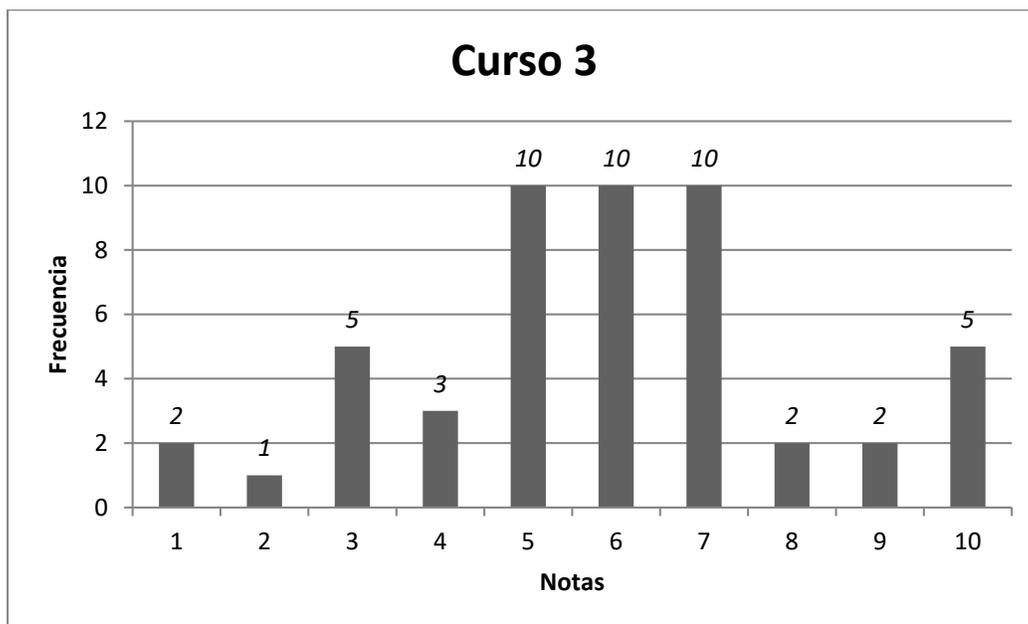


Lo que se observa en el gráfico del resultado del examen de matemáticas del curso 1, es que la mayoría obtuvo una nota de 6 (15 de los 50 alumnos, esto es el 30% de ellos). Siendo el valor modal el que concentra la mayor frecuencia, se tiene que la nota 6 es el modo de esta distribución.



En este gráfico se muestra la distribución de frecuencias del curso 2. Observamos que los alumnos que obtuvieron como nota 4 y 7 se presentan con igual frecuencia, y que estas son las más altas de toda la distribución. Esto significa que el 40% de los alumnos se han sacado notas de 4 y 7 respectivamente, en igual proporción. La

distribución en este caso se dice que es bimodal, siendo los modos la nota 4 y la 7.



La distribución de las frecuencias de las notas del curso 3 se diferencia de las otras en que se observa que las frecuencias más altas se agrupan en torno a las notas 5, 6 y 7. Se trata entonces de una distribución multimodal (es decir, que contiene más de un modo).

Un aspecto destacado del modo como medida de tendencia central es que puede aplicarse a variables que no han sido medidas en escala métrica, dado que se basa en identificar la categoría de la variable que concentra la mayor frecuencia. Veamos un ejemplo: En una librería se realiza un conteo de los libros vendidos en el último mes, discriminando el género literario del mismo. El resultado se expresa en la siguiente tabla de frecuencias:

Género de la publicación	Cantidad de ventas
Novela	198
Ensayo	212
Manual	128
Técnico	406
Autoayuda	449
Infantil	228
Audiolibro	41

Como puede apreciarse en esta tabla, el género literario más vendido es el de autoayuda, cuya frecuencia es de 449. De esto se desprende que la categoría modal de

venta es autoayuda.

Una cuestión que debe tenerse en cuenta cuando se identifica el modo en una distribución de frecuencias es que el modo no es la frecuencia más alta, sino la categoría donde ésta se concentra. En este sentido, y volviendo al ejemplo de la librería, el modo es autoayuda y no 449 que es su frecuencia.

Mediana

La mediana representa el valor de la distribución que deja por encima y por debajo la misma cantidad de casos, otra manera de expresarlo es: la mediana es un valor teórico que permite determinar, sobre el recorrido observado de la variable, el 50% de los casos que quedan en una posición superior e inferior. Para ver como se emplea esta medida, tomemos un conjunto ficticio de datos, que llamaremos notas del parcial de Estadística de la comisión A. El conjunto está compuesto por 19 valores que van entre 2 y 10. Al ordenarlos de menor a mayor encontramos que el valor 5 es el punto en que se puede separar al conjunto de datos en dos mitades iguales. Dicho en otras palabras, el valor 5 permite separar los casos en dos mitades de nueve casos.

2, 2, 2, 2, 4, 4, 4, 5, 5, 5, 5, 6, 10, 10, 10, 10, 10, 10, 10

En este punto es necesario reparar en que el valor de mediana se calcula sobre los valores observados de la variable, es por ello que situados en el valor 5 es posible separar dos conjuntos que contienen al 50% de los casos. En la distribución que hemos presentado el valor 5 coincide con un valor efectivamente observado, sin embargo, si la distribución de casos fuera par, deberíamos encontrar un valor de mediana utilizando los valores empíricos del centro del conjunto de datos. Continuando con el ejemplo, digamos que en una segunda instancia rinden el mismo parcial un conjunto de diez estudiantes cuyas notas son las siguientes:

4, 4, 5, 5, {6, 7}, 8, 8, 9, 10

Aquí se aprecia que la mediana es un valor que debe estar situado entre dos valores dados, y por tanto será el promedio entre ellos: $Mdn = \{6 + 7\} / 2 = 6,5$. Ahora, el valor 6,5 es aquel que separa el conjunto de datos en dos mitades iguales, aunque ese mismo valor no pertenece a la distribución original de datos.

La pregunta que sigue es cómo interpretar la mediana de una distribución, y para ello debemos tener en cuenta varios aspectos. Primero, para una variable dada se pueden obtener n valores teóricos, que en el ejemplo que estamos desarrollando es de 1 a 10. Luego, es necesario registrar si efectivamente se observan todos los valores teóricos, y en este caso se tiene que es posible que ello no ocurra. Si tomamos la

segunda distribución de datos se aprecia que no se registraron los valores 1, 2, 3, es decir que hay una restricción de valores posibles. Cuando esto ocurre los casos se acumulan en determinados valores y la mediana debe ser interpretada como valor teórico, no empírico. Para entender este punto, nos centraremos en la segunda distribución, allí claramente la mediana vale 5, pero si decimos que por encima de ese valor se encuentra la mitad de los casos incurriremos en un error: nótese que por encima de cinco hay 8 casos, lo cual no corresponde con el 50% de la distribución. Entonces, al trabajar la mediana como valor teórico estamos utilizando este estadístico para entender el comportamiento de la variable, no para situar a los individuos dentro de la variable.

Existen otros procedimientos para visualizar la distribución de casos y utilizan medidas de posición entre las que se cuenta la mediana. El más común de ellos es la gráfica de cajas que vimos anteriormente. Por ahora baste mencionar que la mediana es un valor teórico que establece un punto de corte en la variable donde es posible repartir el 50% de los casos por encima y por debajo del mismo.

Media Aritmética o Promedio

La media aritmética es la medida que más se utiliza para resumir información. Puede definirse como el valor equidistante de un conjunto de valores y por esta propiedad sería el dato que mejor representa al conjunto. Veamos un ejemplo de uso del promedio.

En la siguiente tabla se muestra un resumen del resultado de la aplicación de tres métodos diferentes de aprendizaje de la lectura a escolares de primaria. Cada método ha recibido una denominación: A, B y C; la tabla expresa la cantidad de palabras que los niños fueron capaces de leer correctamente.

Método A	Método B	Método C
22	35	44
18	26	49
19	30	51
30	30	52
25	12	39
12	17	43
15	28	50

Una manera de determinar cuál ha sido el método que da mejores resultados, es mediante la comparación del rendimiento promedio de cada uno de los grupos. Matemáticamente el promedio se define como la sumatoria de los valores de la variable, dividido el total de casos. La fórmula conceptual que resume el cálculo es la

siguiente:

$$\bar{x} = \frac{\sum_{i=1}^n x}{n}$$

\bar{x} : media o promedio aritmético

$\sum_{i=1}^n x$: indica la sumatoria de todos los valores de la variable

n : total de casos

El cálculo del promedio para los tres grupos en la prueba de lectura se muestra a continuación:

Método A	Método B	Método C
20,14	25,42	46,85

El promedio estaría reflejando el rendimiento general en lectura de los escolares que aprendieron bajo tres métodos diferentes. Nótese que el promedio es un único valor que debe expresar la centralidad de un conjunto de valores. Por tanto, en la distribución original, algunos de ellos quedarán por encima del promedio y otros quedarán por debajo de él. Por ejemplo, en el caso de la distribución de valores del método A, se tiene que el valor 30 se encuentra por encima del promedio, mientras que el valor 15 se encuentra por debajo del mismo. Adviértase que si el conjunto de datos no hubiera mostrado variabilidad (es decir, todos los escolares hubieran leído la misma cantidad de palabras), no haría falta calcular un promedio.

La variación de los valores en torno al promedio se denomina dispersión, y es posible obtener medidas de esta propiedad de los datos.

Varianza y Desviación Estándar

Veamos gráficamente que se quiere expresar con variaciones en torno al promedio. En el siguiente renglón se han ordenado los valores originales del método A, incluyendo el valor promedio:

12 - 15 - 18 - 19 - 20,14 - 22 - 25 - 30

Es de notar que el promedio es un valor que representa al conjunto de valores, los cuales pueden estar próximos o distantes de él. Mientras más cercanos entre si los valores originales, mayor será la proximidad de los mismos a su promedio y viceversa.

La medida que expresa la distancia del conjunto de valores al promedio, se denomina varianza. Conceptualmente se define como el promedio de las diferencias

cuadráticas entre la media y los valores originales. La fórmula que lo expresa es la siguiente:

$$Var = \frac{\sum_{i=1}^n (x - \bar{x})^2}{n}$$

$\sum_{i=1}^n (x - \bar{x})^2$: sumatoria de las diferencias de cada valor de x respecto a la media
 n = cantidad de casos.

La fórmula de cálculo es similar a la del promedio, solo que en este caso el numerador indica que se debe realizar la sumatoria de las diferencias entre cada valor individual y el promedio. Dado que hay valores que se encuentran por encima del promedio y otros que se encuentran por debajo de él, el resultado debe elevarse al cuadrado. De otro modo, el resultado final sería siempre cero.

Es importante comprender que la varianza está expresando la dispersión de los valores originales respecto del promedio, y por tanto indica que tan homogéneos son esos valores. Dado que esta medida de dispersión esta expresada por valores elevados al cuadrado, se reemplaza por otra medida llamada **Desvío Estándar**. El desvío estándar es la raíz cuadrada de la varianza y su interpretación es equivalente en tanto medida de dispersión, pero con la particularidad de que los valores están en la misma unidad que la media.

$$de = \sqrt{\frac{\sum_{i=1}^n (x - \bar{x})^2}{n}}$$

Veamos cómo se presentan las varianzas y los desvíos estándar de los tres métodos de enseñanza de la lectura

	Método A	Método B	Método C
Media	20,14	25,42	46,85
Varianza	37,14	65,28	23,8
D. E.	6,09	8,07	4,87

Dijimos al principio que necesitábamos una aproximación para determinar cuál método de aprendizaje de la lectura era el más eficiente, o al menos el que producía mejores resultados. Estas medidas en conjunto nos permiten aproximarnos con bastante precisión a tal respuesta:

- a) El método de lectura C, es el que muestra el promedio más alto y la menor varianza, por tanto es el que ha dado los mejores resultados.
- b) El método B produce mejores resultados que el método A, puesto que su promedio es más alto. Pero la dispersión observada es mayor, de modo que el método A produce resultados más bajos en comparación con el método B, pero es más homogéneo en su composición.
- c) El método A es el que produce los resultados más bajos.

La lectura del promedio junto con la desviación estándar nos estaría permitiendo ver que aquellos escolares que aprendieron con el método C, elevaron su rendimiento de manera pareja. En cambio, aquellos que aprendieron con el método B, solo algunos mejoraron su rendimiento, y finalmente los que aprendieron con el método A tuvieron un rendimiento más parejo, pero inferior al obtenido por el uso de los otros métodos de aprendizaje.

En el ejemplo anterior, los datos están disponibles y son solo siete casos, lo cual facilita la lectura de los estadísticos y los casos simultáneamente. Pero, el promedio, la varianza y la desviación estándar, se calculan para resumir conjuntos grandes de datos. Veamos el siguiente ejemplo: en las tablas que se muestra más abajo se presentan los resultados de un muestreo sobre ausentismo en el 3º ciclo de EGB y el Ciclo de Especialización. La variable ha sido medida en dos momentos diferentes sobre las mismas unidades de análisis (alumnos). Para interpretar el conjunto de datos recogido se han resumido en los estadísticos de media, varianza y desviación estándar.

	Primer Cuatrimestre		Segundo Cuatrimestre	
	3º EGB	CE	3º EGB	CE
Media	16	15,83	9,33	16,33
Varianza	56	60,56	13,46	4,667
DE	7,48	7,78	3,66	2,16

Los datos que se presentan muestran que en el primer cuatrimestre, tanto en el 3º ciclo de EGB como en el Ciclo de Especialización, la media y la dispersión en la cantidad de faltas de los alumnos es similar. Sin embargo, en el segundo cuatrimestre es factible observar que el promedio de faltas del 3º ciclo de EGB es más bajo, mientras que en Ciclo de Especialización el promedio se mantiene muy similar. Paralelamente, la varianza del 3º ciclo de EGB es menor en el segundo cuatrimestre, lo cual refleja que en conjunto, los escolares están faltando menos. Para el Ciclo de Especialización, la varianza también es menor, pero dado que el promedio sigue siendo prácticamente muy similar al del primer cuatrimestre, es factible deducir que el grupo tiende a

mostrar un comportamiento particular y es que aquellos escolares que faltaban poco a la escuela, ahora faltan más y aquellos que faltaban mucho, ahora faltan menos.

Conociendo la media y el desvío estándar, se puede determinar la posición relativa de un individuo en relación al grupo. Así por ejemplo, de un alumno que obtiene una calificación de 6 en matemáticas y de 8 en inglés, nos inclinaríamos a pensar que le fue mejor en ésta última materia. Sin embargo si consideramos la media y el desvío estándar del grupo en ambas materias, esta primera impresión resulta injustificada.

Materia	Media	D. E.
Matemáticas	4.0	1.68
Inglés	8.55	1.34

Nótese ahora que una calificación de 6 en matemáticas está por encima de la media en más de una desviación estándar, mientras que una calificación de 8 en inglés está a menos de una desviación estándar por debajo de la media.

Lo antes dicho nos lleva a enfatizar que todas las medidas de tendencia central deben acompañarse de medidas de dispersión, ambas nos dan una aproximación más certera del comportamiento de la variable en el conjunto de unidades de análisis estudiadas.

Coeficiente de variación de Pearson

El coeficiente de variación V de Pearson, viene a resolver el problema de comparar la variabilidad de diferentes conjuntos de datos, que han sido medidos con distintas escalas. Se define como el cociente entre la desviación estándar de la distribución y su media aritmética. Como se deduce, el coeficiente es una manera de determinar cuántas veces la desviación estándar contiene a la media, y por tanto nos informa de la dispersión general. Dado que no depende de la escala en que ha sido medida la variable, lo hace apto para comparar distintas distribuciones. Una medida general del coeficiente se logra multiplicándolo por cien, y en este caso es una medida porcentual de variación. Su formulación matemática es la siguiente:

$$v = \frac{\sigma}{\bar{x}}$$

Veamos el siguiente ejemplo para ilustrar su utilidad: Un psicólogo educacional tiene datos de tres muestras de individuos que han sido evaluados con diferentes pruebas de inteligencia. Desea determinar el nivel de dispersión para cada una de ellas y

verificar cuál es la que tiene menos dispersión. En esta situación la desviación estándar no es la medida óptima pues al usarse pruebas diferentes de inteligencia, no se puede comparar entre sí. Para ello aplica el coeficiente V de Pearson, que para este ejemplo se obtienen los siguientes resultados.

	Media	Desviación Estándar	Coeficiente de Variación (%)
Muestra A	8	5,51	68,87
Muestra B	125	46	36,8
Muestra C	510	212	41,57

Independientemente de la prueba de inteligencia utilizada, se observa que la muestra que menor dispersión evidencia es la B.

Procedimiento para el cálculo del promedio, la varianza y la desviación estándar

Ya no se calculan a mano los estadísticos que hemos repasado hasta aquí, puesto que existen programas con esas funciones, incluso una calculadora científica puede hacer ese trabajo. El procedimiento que mostramos a continuación tiene la finalidad de ilustrar la aplicación de la fórmula de cálculo de la media, la varianza y la desviación estándar. Se basa en la construcción de una tabla que ayuda en la serie de cálculos que estas medidas requieren.

A continuación, se presenta una serie de puntajes obtenidos por siete alumnos en una prueba de ortografía, donde se contabilizó la cantidad de palabras escritas correctamente al dictado.

Alumno	Cantidad de palabras correctas
1	28
2	36
3	40
4	30
5	40
6	25
7	35

Para el cálculo de los estadísticos requeridos, primero es necesario agregar algunas

columnas a la tabla; estas columnas son:

$|x - \bar{x}|$ representa cada valor individual menos la media, sin considerar el signo de la sustracción.

$(x - \bar{x})^2$ elevar al cuadrado el resultado de la sustracción anterior.

Alumno	Cantidad de palabras correctas	$ x - \bar{x} $	$(x - \bar{x})^2$
1	28	5,42	29,37
2	36	2,58	6,65
3	40	6,58	43,29
4	30	3,42	11,69
5	40	6,58	43,29
6	25	8,42	70,89
7	35	1,58	2,49
n=7	$\Sigma=234$		$\Sigma=205,18$

Con las dos primeras columnas de la tabla y aplicando la fórmula presentada anteriormente, se obtiene el valor de la media:

$$\bar{x} = 234/7 = 33,42.$$

Lo que sigue es restar este valor a cada observación para completar la siguiente columna. De este modo, para el alumno 1 se tendrá: $28 - 33,42 = -5,42$, puesto que no consideramos el signo de la resta utilizaremos el valor absoluto: 5,42. La siguiente columna se completa elevando al cuadrado el resultado obtenido: $(5,42)^2 = 29,37$. Realizando esta operación para cada uno de los alumnos, obtendremos la columna final, cuya sumatoria nos permitirá obtener la varianza:

$$\text{Var} = 205,18/7 = 29,31.$$

La desviación estándar, es la raíz cuadrada de la varianza, por lo tanto:

$$\text{DE} = \sqrt{29,31} = 5,41.$$

Así, para el conjunto de datos analizados se tiene que

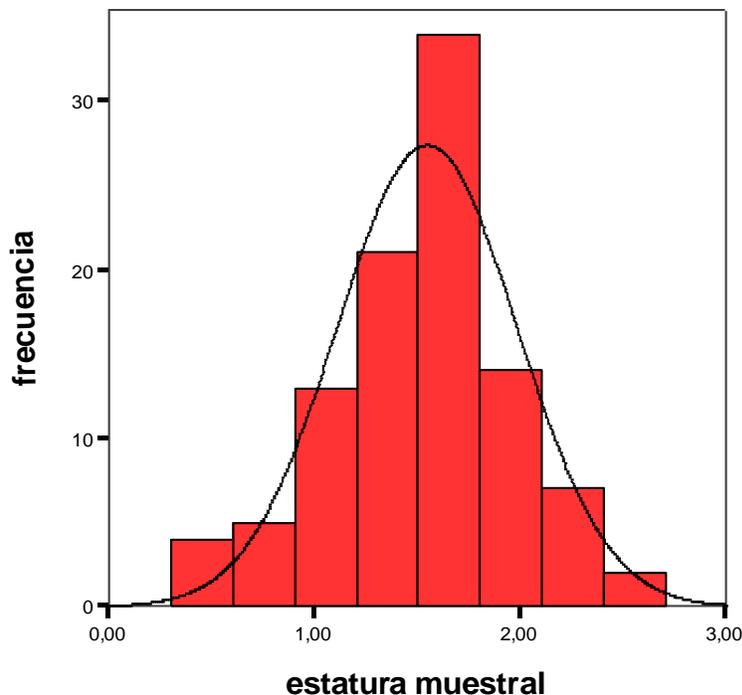
Promedio	Varianza	Desviación Estándar
33,42	29,31	5,41

Capítulo 5

La distribución normal

En nuestra vida cotidiana estamos utilizando casi todo el tiempo el concepto de normalidad estadística, aunque no somos conscientes de ello. Tal concepto nos permite hacer juicios de valor y establecer generalidades acerca de aspectos de la vida que son los más frecuentes de observar, y nos permite también comprender fácilmente a qué nos referimos cuando valoramos un evento o acontecimiento como muy poco probable.

Vamos a tomar una variable aleatoria en la población, tal como es la estatura de las personas. Decimos en este caso que la variable es aleatoria porque en cada medición que realicemos, el valor puede variar libremente. Sin embargo, la mayoría de las medidas tomadas en la población, tenderán a centrarse en un determinado valor, el cual aparecerá como el más frecuente. Para este ejemplo, tomamos datos de salud pública del Ministerio de Salud de la Nación del año 2010. Seleccionamos 100 personas con edades entre cinco meses y veintiún años de edad. La selección de la muestra es enteramente accidental. Con estos datos construimos la siguiente gráfica.



Además, calculamos algunos estadísticos sobre la muestra que se presentan en la siguiente tabla.

Estadístico de la variable estatura muestral

Media	Mediana	Modo	Desviación estándar	Varianza	Máximo	Mínimo
1,54	1,64	1,70	0,44	0,19	2,53	0,5

*Los valores están expresados en metros

En este ejemplo el valor modal de la distribución es 1,70 mts. que está representado por el punto más alto de la distribución, luego vemos que las estaturas bajas son menos frecuentes y lo mismo que las estaturas altas. El centro de la distribución está próximo al modo ya que el promedio es de 1,54 mts. para esta muestra. Basándonos en la información que tenemos y en el gráfico diremos que aquellos individuos por debajo de 1,7 mts. deberían considerarse bajos, mientras que aquellos por encima de esa estatura se consideran individuos altos. En la tabla de frecuencia original, aparecen tres casos con una estatura de 0,60 mts. los cuales los consideramos demasiado bajos. También aparece un caso con una estatura de 2,53 mts., al cual lo consideramos como demasiado alto. Estadísticamente diremos que estos casos son atípicos o extremos dentro de la distribución que estamos analizando.

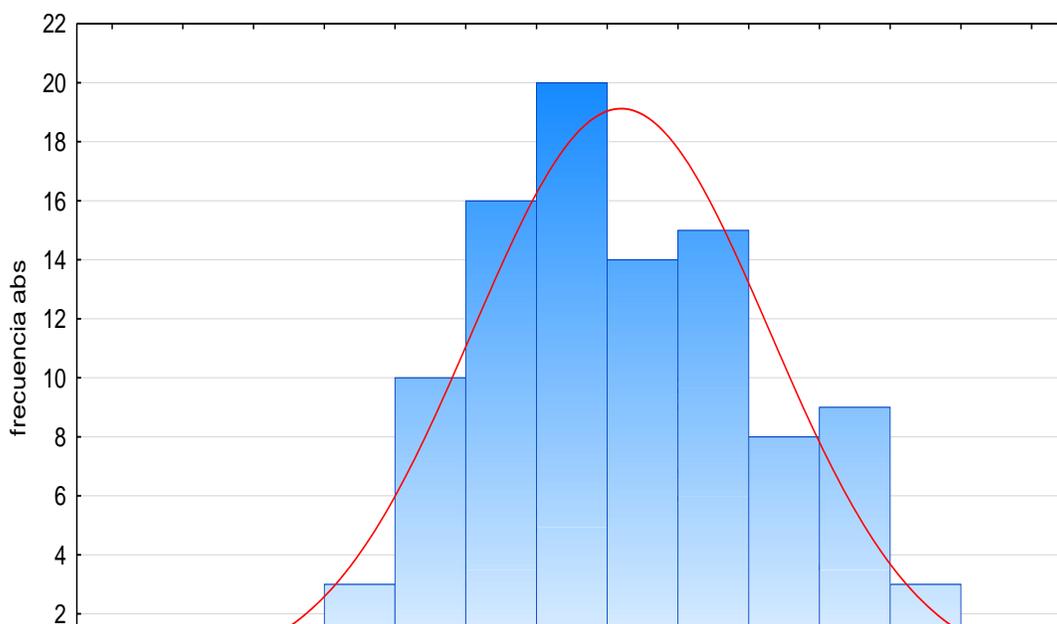
Sobre el histograma hemos superpuesto una curva con forma de campana, conocida como distribución normal. Esa distribución corresponde a una curva teórica en donde las estaturas se acumulan proporcionalmente en torno a una media de 1,54 mts. y una desviación estándar de 0,44. Vemos que la superposición entre la curva normal y el histograma coinciden en varios puntos, por lo cual es factible utilizar el modelo de distribución normal para la variable estatura en esta muestra.

Resumiendo lo dicho hasta aquí tenemos que lo que nos informaría esta curva para la variable estatura es que la mayoría de las personas están centradas en el valor 1,70 mts, luego los más altos se acumularían hacia la derecha de la distribución y los más bajos hacia la izquierda. Como los individuos altos y muy altos (como así también los bajos y muy bajos), no son comunes en la población, tenemos que su frecuencia disminuye a medida que nos alejamos del valor central, lo que confiere a la distribución su típica forma acampanada. Debemos tener siempre presente que la altura de la curva indica frecuencia.

Volvamos ahora sobre los casos atípicos o extremos: ¿cómo explicar que en una población existan casos de individuos con 60 centímetros de estatura y otro con 2,53 metros? Recordemos que la muestra es accidental y obtenida de un listado de salud pública del Ministerio de Salud de la Nación del año 2010. Ese listado incluye persona de todas las edades y al cotejar la edad de los individuos de 60 cm. de estatura

encontramos que son menores de un año de edad. Por lo tanto, no pueden considerarse como casos atípicos. ¿Qué hay de una persona de 2,53 mts? En el listado original no había ninguna persona de esa estatura, con la intención de mostrar un caso verdaderamente atípico incluimos la estatura de Trijnje Kever, quien se supone fue la mujer más alta de los Países Bajos, y según se reporta vivió solo 17 años. En el listado original la estatura promedio de quienes tienen 17 y 18 años es de 1,69 mts. por lo cual encontrar una mujer de 17 años de 2,53 sí puede considerarse como un caso excepcional o atípico.

La distribución normal tiene importantes propiedades para la estadística descriptiva e inferencial, dado que puede utilizarse como modelo matemático para analizar variables aleatorias relevantes para la investigación educacional. El siguiente gráfico puede ilustrar lo dicho dado que se trata de un conjunto de valores generados aleatoriamente tomando como referencia una media igual a 100 y varianza igual a 500. Lo interesante de esta simulación es que el patrón de valores individuales fue generado en función de los resultados de la prueba PISA aplicada en 2012 para matemáticas.



Nótese que este simple ejercicio genera un patrón de datos que puede ser analizado bajo los postulados de la distribución normal estandarizada. En el apartado que sigue nos ocuparemos de esta propiedad de la distribución normal.

La distribución normal como modelo matemático: la normal estandarizada

Antes de comenzar este apartado es importante reconocer que las propiedades del modelo matemático implícito en la distribución normal estandarizada, son aplicables a una distribución empírica solo cuando ésta última se aproxima a la normal. Aunque no es objeto de desarrollo aquí, existen procedimientos de transformación de las puntuaciones originales de una distribución empírica para que se empareje con una distribución normal y así utilizar sus propiedades. Un ejemplo de esto lo vimos anteriormente cuando simulamos una muestra aleatoria de las puntuaciones PISA con una media de 100 y varianza de 500.

La distribución normal es unimodal y simétrica, lo cual significa que la media, la mediana y el modo coinciden en el mismo valor, que es el centro de la distribución. La simetría indica que la distribución es exactamente igual hacia ambos lados de la media.

La distribución normal es continua y asintótica (esto es, que se extiende indefinidamente en sus extremos), esto quiere decir que cualquier valor de la variable estaría contenido dentro de la distribución normal estandarizada. Las variables empíricas, como sabemos, no contienen infinitos valores, por lo cual, dentro de ciertos límites una distribución normal contendrá el 100% de los valores de esa distribución empírica. Esto quiere decir que la varianza de una distribución empírica es finita; otra manera de expresarlo es diciendo que todo conjunto de datos tiene un valor máximo y un mínimo. El teorema de Chebyshev proporciona una regla empírica que determina que dado un conjunto de valores que se asemeje a una distribución normal, aproximadamente el 90% de los valores de esa distribución estará comprendida entre dos desviaciones estándar por encima y por debajo de la media. Los detalles de este teorema están desarrollados en la desigualdad de Chebyshev, pero no lo desarrollaremos en este tutorial.

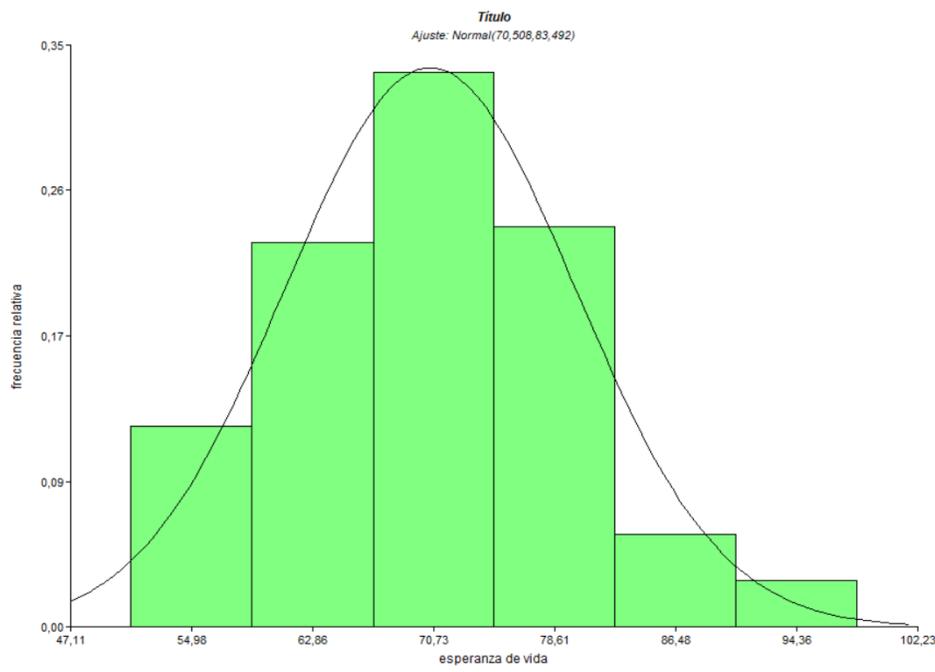
Los valores de media y desviación estándar en una distribución normal estandarizada están expresados en puntuación Z. De esto se deriva que la media de una distribución normal estándar será siempre cero y la desviación estándar será siempre uno. Esta es una propiedad muy importante dado que cualquier distribución empírica puede ser estandarizada, transformando sus puntuaciones originales en puntuación Z aplicando la siguiente fórmula:

$$Z = \frac{x - \bar{x}}{DE}$$

Veamos como ejemplo la siguiente tabla extraída del ministerio de Salud de la Nación del año 2010, en la que se obtuvo una muestra aleatoria de 109 casos donde se registró la esperanza de vida. La distribución obtenida se aproxima a la distribución normal, por lo tanto, es lícito realizar la transformación Z.

Esperanza de vida en la población

Casos	media	DE	mínimo	máximo
109	70.15	10.57	43	82



Como se aprecia en el gráfico, la distribución empírica se adapta a una distribución normal estándar. Del total de casos analizado elegimos tres y transformamos sus puntajes a puntuación Z, tal como lo muestra la siguiente tabla

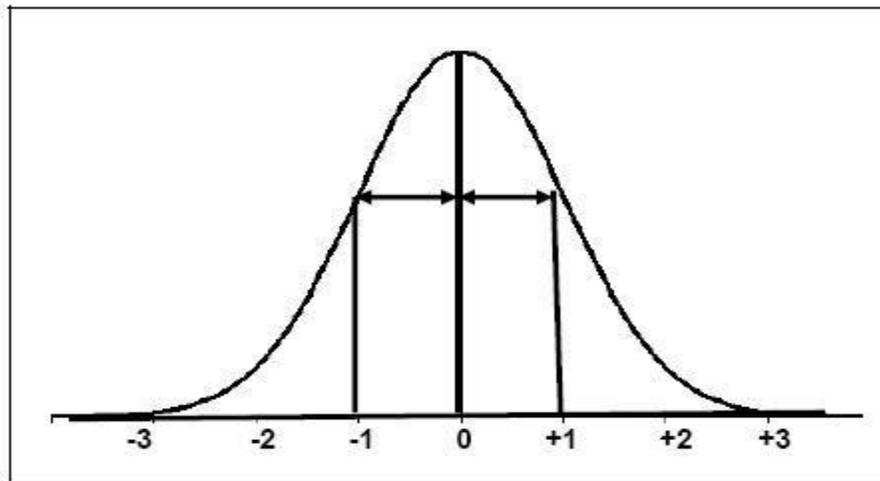
Puntaje original	Puntaje Z
79,00	0,84
80,00	0,93
53,00	-1,62

Los puntajes representados en la tabla son equivalentes, la diferencia radica en que el

puntaje original está centrado en la media de 70,15 y la transformación Z arroja una media de 0. Veamos ahora cómo queda la tabla si calculamos los mismos estadísticos sobre la distribución esperanza de vida a valores Z.

Casos	media	DE	mínimo	máximo
109	70.15	10.57	43	82
109	0.00	1.00	-2.57	1.12

De esta forma la distribución de puntuaciones originales transformadas de una distribución empírica, queda asimilada en la siguiente distribución, cualquiera sea la puntuación original. Como vimos, la única condición para esto es que la variable empírica se aproxime con fidelidad a una distribución normal teórica.



Volvamos un instante al ejemplo de la distribución de las estaturas. Haciendo un recorte de la muestra de estaturas para las personas que tiene 17 y 18 años, tenemos que la media es de 1,69 mts. y la desviación estándar es de 0,22 mts. ¿Qué tan diferente resulta la estatura de Trijnje Keever? Si observamos el gráfico anterior vemos que la mayoría de las personas se encuentran entre ± 1 desviación estándar de la media, más infrecuentes son las que se apartan ± 2 desviaciones estándar y aquellas a ± 3 desviaciones estándar son sin duda casos atípicos, o muy raros. La estatura de Trijnje Keever se aparta 3.81 desviaciones estándar de la media.

$$Z = \frac{2,53 - 1,69}{0,22} = 3.81$$

Trijnje Keever murió a los 17 años en 1633 y hasta el presente no se ha registrado una altura igual en una mujer.

Sobre la distribución estandarizada podemos trabajar a la inversa. Es decir, conociendo los valores en puntuación Z, podemos descomponer el valor original. Recordemos del ejemplo de la esperanza de vida de la población que su media era de 70,15 años y su desviación estándar de 10,57 años. Para obtener los puntajes originales tendríamos que emplear la siguiente ecuación:

$$Po = \bar{x} + (Z * de)$$

Siendo Po la puntuación original, la tabla a continuación muestra el cálculo de transformación de puntaje Z a puntuación original.

Puntaje Z		Po
0,84	$70,15+(0,84*10,57)$	79,0
0,93	$70,15+(0,93*10,57)$	80,0
-1,62	$70,15+(-1,62*10,57)$	53,0

La distribución normal estandarizada y la proporción de casos

Hasta aquí hemos podido ver que mediante la estandarización de una distribución empírica, hacemos razonable y comprensible los términos de normal o típico, o bien, extremo o atípico. También podemos trabajar la proporción de casos que son típicos y atípicos, mediante la utilización de las áreas bajo la curva normal en término de proporciones.

La siguiente figura nos muestra en el eje de las abscisas las puntuaciones estandarizadas Z, las cuales se proyectan sobre la curva y en el punto donde se intersectan, se define un área bajo la curva. De este modo, si tomamos el segmento comprendido entre los valores $Z = \pm 1.04$, tenemos que se ha cubierto el 70% del área bajo la curva. Esto tiene una utilidad fundamental para comprender el comportamiento de variables empíricas en poblaciones e intentaremos mostrarlo mediante un ejemplo. Supongamos que se ha tomado una prueba de rendimiento en lectura en la provincia de Córdoba, cuya distribución de valores se aproxima a una distribución normal. La media en dicha prueba es de 550, y la desviación estándar es de 125. La curva que se muestra en la figura que sigue podría utilizarse como modelo dimensional de esa variable y mediante algunos simples cálculos podríamos contestar

las siguientes preguntas:

1 ¿Qué puntuaciones originales de la prueba comprenden el 70% de la población, considerada como rendimiento medio?

2 ¿Qué puntuación le corresponde al 5% de la población con el rendimiento más bajo?

3 ¿Qué puntuación puede considerarse excepcionalmente alta en la prueba?

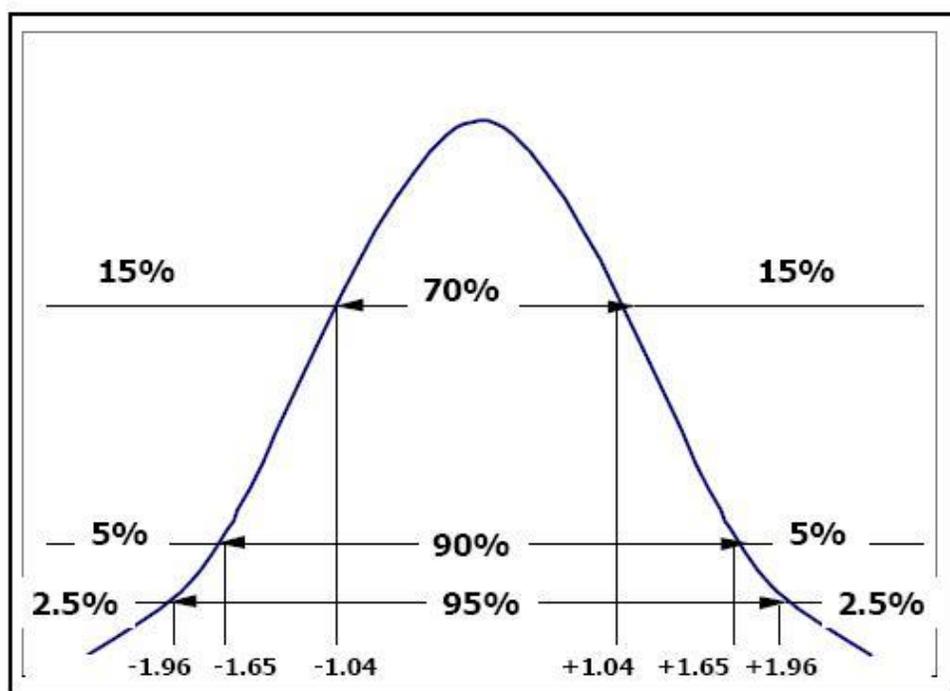


Figura publicada en: Morales Vallejo, Pedro (2008) Estadística aplicada a las Ciencias Sociales. Madrid: Universidad Pontificia Comillas

Para responder a la primera pregunta debemos aplicar la fórmula mediante la cual transformamos la puntuación Z en puntuación original. Los valores Z que nos informan qué valores contienen al 70% de la población con puntuación en los valores medios de la distribución son, como vimos $\pm 1,04$. Pasándolos a puntuación original tenemos que el 70% de la población con puntuación en los valores medios de la prueba han obtenido puntajes de 420 como mínimo y 680 como máximo.

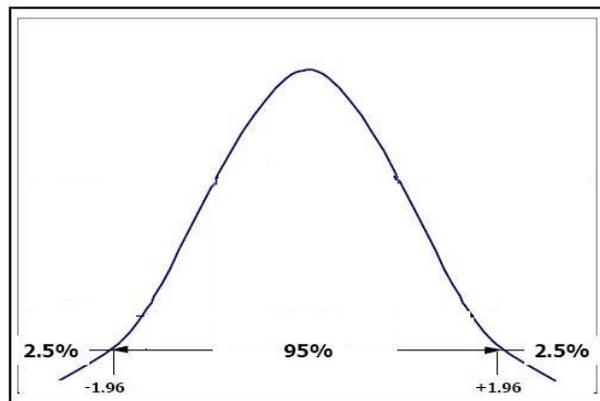
Para responder a la segunda pregunta debemos considerar el valor Z que deja por debajo el 5% del total de valores contenidos en la distribución. Dicho de otro modo, debemos buscar el valor correspondiente al punto de corte para el 5% hacia la izquierda de la distribución; ese valor es -1.65 . Nuevamente, transformando ese valor

a puntuación original tenemos que la puntuación original de prueba obtenemos el valor $343,75 \approx 344$ (por redondeo). Por lo tanto todos aquellos alumnos que hubieran obtenido esa puntuación o menor, estarían en el 5% de la población con más bajo rendimiento.

Por último, una puntuación excepcionalmente alta podría considerarse aquella que corresponde al 1% de la población con más alto rendimiento. Repitiendo el procedimiento, tenemos que el valor Z de 2,34 deja por debajo de al 99,04% de la distribución. Es decir que por encima de ese valor se encuentra el 1% de la población aproximadamente (este valor no consta en la figura, se extrajo de una tabla que describe los valores Z y el porcentaje de área correspondiente). Transformando ese valor Z a puntuación original tenemos que aquellos individuos con un puntaje de 842,5 o más en la prueba de lectura, serían los que ocuparían el 1% de la distribución con rendimiento excepcionalmente alto.

Áreas bajo la curva como distribución de probabilidades

Más adelante veremos cómo la distribución de probabilidades sirve a distintos propósitos, especialmente para la prueba de hipótesis, por esto merece una referencia el porcentaje de área bajo la distribución normal como distribución de probabilidades. Para ello vamos a prestar especial atención al valor Z de la curva $\pm 1,96$. Veamos la siguiente figura:



Vemos que los valores $\pm 1,96$ contienen el 95% de casos de la distribución, y dejan fuera el 5% de los restantes casos. Ese 5% se reparte entre aquellos que ocupan ambos extremos de la distribución, por tanto se cuentan 2,5% en la parte izquierda (aquellos que están por debajo del valor Z -1,96), y un 2,5% en la parte derecha

(aquellos que están por encima del valor $Z + 1,96$).

Si interpretamos esos porcentajes en término de probabilidades, y tenemos que escoger un valor al azar del total de valores de la variable en esa distribución, es muy probable que el valor escogido sea mayor o igual que $-1,96$ y menor o igual que $+ 1,96$. Deducimos esto porque el 95% de los valores de la distribución se halla contenido entre esos dos puntajes Z . Esto equivale a decir que obtener por azar un valor menor que $-1,96$, o mayor que $+ 1,96$ es muy improbable. Esto se debe a que por debajo y por encima de esos valores solo se halla contenida el 5% de la distribución.

Veamos con un sencillo ejemplo cómo esta propiedad de la distribución normal puede ayudarnos a esclarecer algunos aspectos en las conclusiones o afirmaciones de un estudio (más adelante veremos que es el procedimiento usual en la prueba de hipótesis). Supongamos que la evaluación provincial de lectura, arroja una media de 550. Supongamos ahora que la media de Córdoba Capital en esa evaluación es de 643 con una desviación estándar de 104. Al ver la diferencia del promedio provincial con el de Capital, podríamos afirmar que el rendimiento lector de los escolares cordobeses que residen en la capital es significativamente diferente del resto de la provincia. Entendamos el término significativamente como una diferencia sistemática suficientemente amplia que debe interpretarse como que los escolares residentes en Capital leen mejor que el resto.

Podemos utilizar la distribución normal para testear si esa afirmación puede sostenerse empíricamente o solo debe tomarse como una afirmación basada en una simple diferencia. Dijimos que la distribución de las puntuaciones de esa evaluación seguía un modelo normal, por lo cual dichas puntuaciones pueden transformarse a puntaje Z . Si los escolares de Capital tienen un rendimiento significativamente superior al resto, la puntuación Z de esa localidad debería aparecer en la distribución como un valor atípico. Dicho en otras palabras, su puntaje Z debería corresponder a la zona delimitada por el 2,5% superior de la distribución. Como ya sabemos, ese puntaje es $Z=1,96$.

Nos resta transformar la puntuación original de Capital a valor Z y comprobar donde queda situada en la distribución. Para esto empleamos la fórmula ya conocida y tenemos que:

$$Z_{capital} = \frac{643 - 550}{104} = 0,89$$

Como se aprecia el valor Z de Capital en la distribución es menor que el valor crítico 1,96 sobre el cual estarían los valores correspondientes al 2,5% de mayor rendimiento. Es factible por tanto concluir que, aunque el promedio de Capital es

mayor que el del resto de la provincia, tal diferencia no es suficiente para afirmar que sea significativa.

La posibilidad de utilizar la distribución normal como distribución de probabilidades, resulta muy útil en el estudio de poblaciones, especialmente cuando nuestro foco de atención son variables que se aproximan a esa distribución teórica. Puesto que la distribución de cualquier variable con distribución normal puede ser estandarizada; esto es, transformada a valor Z, es factible calcular la probabilidad de que un valor en la variable sea mayor o menor a un valor Z dado.

Existen maneras de calcular los valores Z mediante las tablas de distribución normal que se publican en los textos de estadística. Aquí recurriremos a una función de un programa estadístico gratuito que se encuentra en el siguiente enlace.

<http://stattrek.com/online-calculator/normal.aspx>

Al ingresar al sitio encontraremos una pantalla como la siguiente:

Normal Distribution Calculator: Online Statistical Table

The Normal Distribution Calculator makes it easy to compute cumulative probability, given a normal random variable; and vice versa. For help in using the calculator, read the [Frequently-Asked Questions](#) or review the [Sample Problems](#).

To learn more about the normal distribution, go to Stat Trek's [tutorial on the normal distribution](#).

■ Enter a value in three of the four text boxes.
■ Leave the fourth text box blank.
■ Click the Calculate button to compute a value for the blank text box.

Standard score (z)

Cumulative probability $P(Z \leq z)$

Mean

Standard deviation

Calculate

Aunque la página está en inglés, nos interesa solo los comandos que aparecen señalados como:

Standard score (z): Puntaje estandarizado Z.

Cumulative Probability $P(Z \leq z)$: (Probabilidad acumulada hasta un valor Z dado).

Mean=0: (Media= 0)

Standard deviation=1: (desviación estándar= 1)

La función puntaje estandarizado Z, nos permite ingresar un valor Z conocido, y nos devuelve en la función Probabilidad acumulada hasta ese valor Z. Las restantes funciones indican que estamos trabajando sobre una distribución normal estandarizada con media igual a cero y desviación estándar igual a uno.

Retomemos ahora el ejemplo de la prueba de rendimiento en lectura en la provincia de Córdoba, cuya distribución de valores se aproxima a una distribución normal. Habíamos dicho que la media en dicha prueba es de 550, y la desviación estándar es de 125. En base a estos datos vamos a plantear dos preguntas:

1. ¿Cuál es la probabilidad seleccionar al azar un individuo de esa población, cuya puntuación en la prueba de lectura sea de 650 o más puntos que la media?

Para resolver esta cuestión debemos transformar el puntaje original en puntuación Z mediante la siguiente ecuación:

$$Z = \frac{650 - 550}{125} = 0,8$$

Al ingresar ese valor Z en el cuadro de diálogo, nos da como resultado que la probabilidad acumulada hasta Z=0,8 es de 0,78814, tal como indica la figura siguiente.

■ Enter a value in three of the four text boxes.
■ Leave the fourth text box blank.
■ Click the Calculate button to compute a value for the blank text box.

Standard score (z) 0.8

Cumulative probability: P(Z ≤ 0.8) 0.78814

Mean 0

Standard deviation 1

Calculate

Dado que la pregunta es sobre la probabilidad de hallar un valor igual o mayor que 650, debemos restar ese valor a 1, que sería la probabilidad acumulada total. Por lo tanto tendremos que:

$$P_{z \leq 650} = 1 - 0,78814 = 0,21186$$

Dado que las probabilidades pueden interpretarse como porcentaje, estamos en condiciones de afirmar que aproximadamente el 21% de la población obtendrá un rendimiento en la prueba de lectura de 650 puntos o más.

Esta propiedad de la distribución normal nos permite realizar evaluaciones a grandes conjuntos de la población y luego tipificar sus valores en términos porcentuales. Las pruebas estandarizadas de rendimiento más conocidas, muestran sus valores poblacionales en percentiles, que es una escala Z transformada a valores que van de cero a 100. Esta es una herramienta muy útil para establecer comparaciones entre diversas poblaciones utilizando solo un instrumento de evaluación.

Otra pregunta que podríamos plantearnos es la siguiente, dados dos puntajes de prueba:

2. ¿Qué porcentaje de la población obtiene puntajes en la prueba, comprendidos entre 450 y 650 puntos?

Esta pregunta se plantea en el contexto de este ejemplo como una derivación de lo que hemos esbozado en la respuesta anterior. Nótese que estamos intentando establecer un porcentaje de población entre dos valores dados que son cien puntos por debajo y por encima de la media poblacional. Aquí no podríamos utilizar simplemente la desviación estándar dado que es de 125. En este caso procedemos a calcular los valores Z correspondientes a las puntuaciones originales. Ya sabemos que la puntuación Z para el puntaje 650 es de 0,8; deducimos entonces que el valor Z correspondiente a un valor de 450 es de -0,8. La deducción se basa en que ambos puntajes se apartan de la media en 100 puntos, y en el hecho de que una distribución normal es simétrica.

Asimismo, ya sabemos que el valor de probabilidad acumulada hasta el puntaje $Z=0,8$ es de 0,78814. Nos resta saber cuál es la probabilidad acumulada hasta el valor $Z=-0,8$. Utilizando nuevamente el programa tenemos que ese valor es igual a 0,21186. Como ya lo habrá notado ese es valor de probabilidad por encima de $Z=0,8$, solo que ahora estamos trabajando sobre la otra mitad de la curva. Podríamos haber llegado a ese valor por deducción, pero interesa subrayar el carácter simétrico de la distribución y la aplicación de esta propiedad.

- Enter a value in three of the four text boxes.
- Leave the fourth text box blank.
- Click the Calculate button to compute a value for the blank text box.

Standard score (z)	-0.8
Cumulative probability: P(Z ≤ -0.8)	0.21186
Mean	0
Standard deviation	1

Dado que el valor $p=0,78814$ es la probabilidad acumulada hasta el puntaje $Z=0,8$, debemos ahora sustraer el valor $p=0,21186$, que es la probabilidad acumulada hasta el valor $Z=-0,8$. Esta resta la hacemos a los fines de centrarnos en el porcentaje comprendido entre esos valores; así tenemos que la resta de ambas probabilidades nos da $0,57628$, lo cual transformado a porcentaje nos brinda la respuesta buscada. Aproximadamente el $57,6\%$ de la población, obtendrá puntajes en la prueba de rendimiento lector comprendida entre 450 y 650 puntos.

Una de las propiedades de la distribución normal y la transformación Z es que se pueden hacer comparables dos variables que no están en la misma unidad de medida. Dicho en otros términos, la puntuación Z es una medida estandarizada independiente de la unidad de medición.

En el apartado de medidas de tendencia central, nos encontramos un problema en el cual un alumno obtuvo un 6 en matemáticas y un 8 en inglés. Utilizando la media y el desvío estándar concluimos que la nota 6 era una mejor puntuación en términos relativos. Si las mediciones provinieran de una distribución normal de datos, hubiera sido posible transformar los puntajes de matemática e inglés a valores Z mediante la siguiente fórmula:

$$Z = \frac{(x - \bar{x})}{DE}$$

Y así obtendríamos que:

$$Z_{matemáticas} = \frac{(6 - 4)}{1,68} = 1,19$$

$$Z_{\text{inglés}} = \frac{(8 - 8,55)}{1,34} = -0,41$$

Ahora vemos cómo el rendimiento en matemáticas es mayor que el rendimiento en inglés, mediante la comparación de sus respectivos puntajes Z: $-0,41 < 1,19$.

Las puntuaciones Z, al ser medidas estandarizadas, permiten comparaciones aun cuando la variable hubiera sido medida en la misma escala. Esto facilita la interpretación de la posición relativa de cada unidad de análisis. Por ejemplo, si quisiéramos ubicar la posición relativa de cuatro provincias respecto al promedio nacional en las evaluaciones nacionales sobre rendimiento en matemáticas, se podría utilizar el puntaje Z de la siguiente manera.

	Media	DE
Total País	59	25

Teniendo la media y la desviación estándar para el total del país, es posible posicionar cada provincia mediante su puntaje Z; así para Córdoba esta puntuación será:

$$Z_{\text{Córdoba}} = \frac{(61 - 59)}{25} = 0.08$$

Si procedemos de la misma manera con cada provincia, obtenemos la siguiente tabla:

Provincia	Media	Puntaje Z
Córdoba	61	0.08
Rio Negro	68	0.36
Buenos Aires	63	0.16
Formosa	41	-0.72

Entonces, mediante la puntuación Z es posible establecer cuáles provincias han obtenidos buenos resultados en la evaluación, y cuáles han obtenido una puntuación menos favorable.

Misceláneas

1. La distribución normal y sus propiedades han sido objeto de estudio de matemáticos que han dejado una profunda huella en la historia, tal el caso de Abraham De Moivre, Pierre Simon Laplace y Carl Friedrich Gauss. Sin embargo, el uso y la aplicación extendida al estudio de fenómenos sociales se lo debemos a Lambert

Quetelet y Francis Galton. Pero fue finalmente Karl Pearson quien popularizó el término de curva normal.

2. La curva de distribución normal es un modelo matemático aplicado a la distribución de variables empíricas. La función de densidad de la distribución está dada por la siguiente expresión.

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

Parece una fórmula muy compleja, pero contiene solo dos parámetros que son variables: μ (mu) y σ (sigma), donde:

μ : media de la distribución poblacional

σ : desviación estándar de la distribución poblacional

3. Cuando tratamos con la distribución normal estándar, tenemos que $\mu=0$ y $\sigma=1$. En estadística matemática suele encontrarse que la distribución normal estándar se expresa como: $X \sim N(0,1)$.

4. La utilización de la distribución de las áreas bajo la curva normal que hemos visto en apartados anteriores se lo debemos principalmente al trabajo de Pafnuti Chebyshev, quien propuso su Teorema donde,

Si $X \sim N(\mu, \sigma)$, entonces:

a) $p(\mu-\sigma < X < \mu+\sigma) = 0,68$

b) $p(\mu-2\sigma < X < \mu+2\sigma) = 0,95$

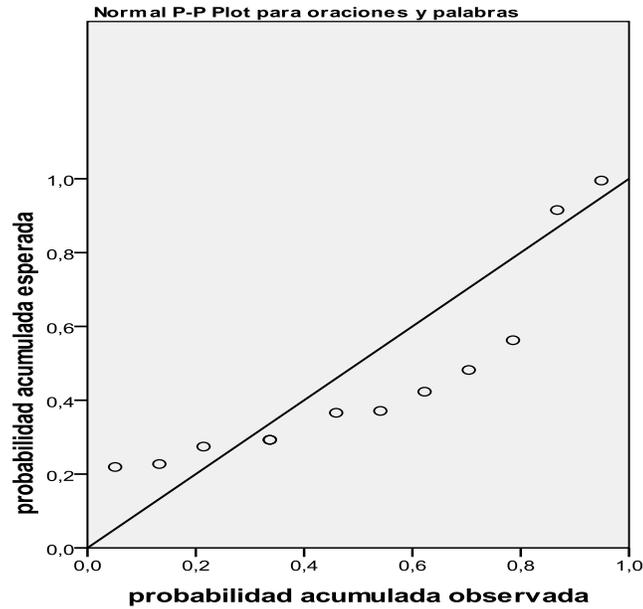
c) $p(\mu-3\sigma < X < \mu+3\sigma) = 0,997$

Lo que equivale a decir que el 68% (aproximadamente) de los valores que tome la variable x estarán situados a una distancia de la media inferior a una desviación estándar. De la misma manera, el 95% de los valores estarán situados a menos de 2 veces la desviación estándar, y un 99,7% de dichos valores se encontrarán dentro de 3 desviaciones estándar. Por lo tanto, para una variable empírica que se aproxime a la distribución normal, la mayor parte de los valores quedan comprendidos a tres desviaciones estándar de la media. Este teorema se volvió fundamental en la investigación y la prueba de hipótesis.

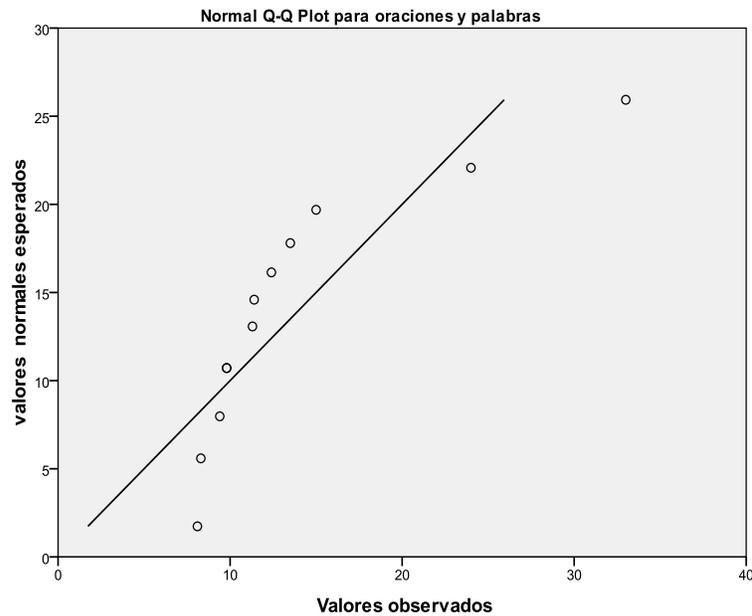
Distribuciones asimétricas

La distribución normal como modelo es una herramienta poderosa para el análisis de datos, pero resulta esencial verificar que una variable empírica se comporta como tal, dado que no podríamos aplicar ninguna de las técnicas estadísticas descriptas si ese no fuera el caso. Las más simples y útiles de las verificaciones son los gráficos de la distribución de los valores de la variable empírica, tales como el histograma o el diagrama de cajas. En ellos, una inspección visual nos informa si la variable sobre la que estamos trabajando se aproxima o no a una distribución normal. Aquí es importante tener en cuenta el término aproximación dado que nunca una variable empírica resultará exactamente una curva normal. Por lo tanto, existen medidas, llamadas coeficientes de asimetría y curtosis que nos dicen qué tan diferente es una distribución empírica de una normal. Solo por citar uno, el coeficiente de asimetría y de curtosis de Fisher nos informa dentro de qué límites podremos considerar a una distribución empírica como una distribución normal.

En los programas estadísticos existen gráficos y pruebas de normalidad que nos permiten verificar rápidamente si una distribución empírica puede tratarse como una normal. Un ejemplo son los gráficos P-P, donde se representan en ejes cartesianos las proporciones acumuladas de la variable empírica, con las de una distribución normal teórica tomada de esos mismos valores. En el gráfico aparece una variable real que se denomina oraciones y palabras. La misma evalúa la capacidad de retención mnémica. Es una variable métrica y se ha aplicado a una muestra de alumnos de primaria. Para saber si la misma sigue una distribución normal bastaría ver su histograma, pero el gráfico P-P estandariza la variable para que quede comprendida entre los valores 0 y 1, luego se grafica la proporción de casos acumulados si la variable tuviera una distribución normal perfecta. En el gráfico se representa como una recta y es la probabilidad esperada ideal para la distribución normal. Entonces, si la distribución empírica está próxima a la normal los casos efectivamente observados estarían muy próximos a la recta (los puntos deberían parearse con la distribución teórica). En el diagrama que vemos más abajo, se aprecia que los casos se apartan de la recta, y por tanto estamos en presencia de una distribución que no se aproxima de manera fiable a una distribución normal. El gráfico recibe el nombre P-P en tanto produce una transformación de probabilidad acumulada basada en el modelo normal.



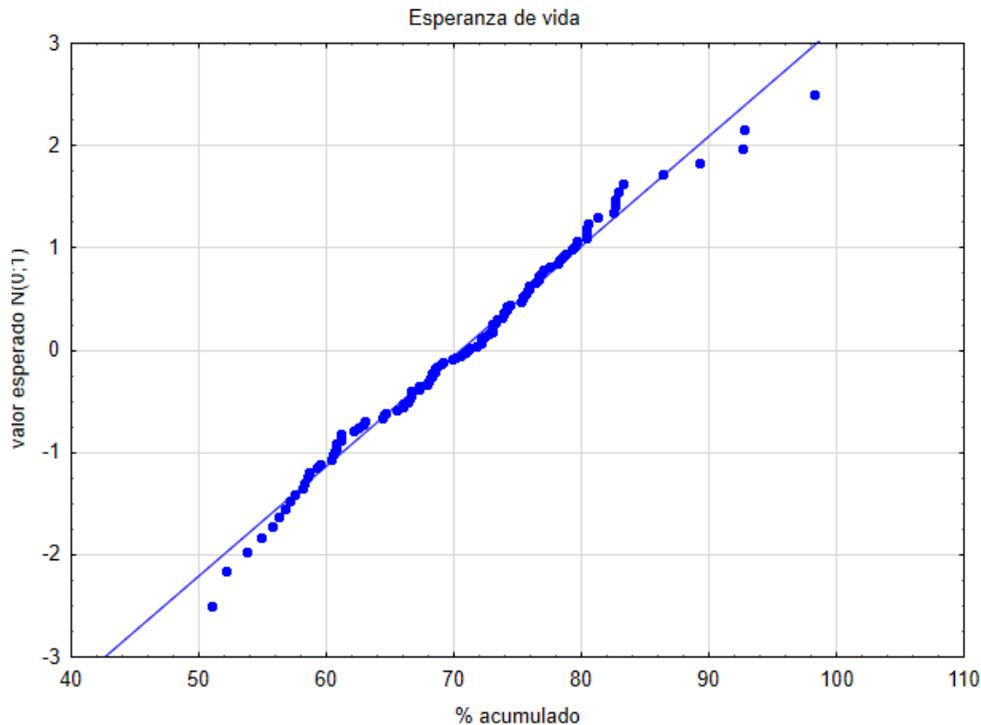
Otros gráficos que siguen la misma lógica son conocidos como gráficos Q-Q donde se representan los cuantiles observados en la distribución empírica respecto a los cuantiles que se esperaría observar si los datos provinieran de una distribución normal. Los cuantiles son una estandarización, pero que se toma en base 100.



La interpretación de este diagrama es la misma que para los gráficos P-P, pero basada en los percentiles.

Además de permitir valorar la desviación de la normalidad, estos gráficos permiten determinar si las curvas observadas muestran desviaciones atendibles de la simetría (curvas en forma de U), o la curtosis (curvas en forma de S).

Los ejemplos que vimos anteriormente representan curvas que se apartan de la normalidad estadística. Pero, hemos trabajado con una variable llamada esperanza de vida en la población que tiene una gran proximidad con la distribución normal. En una gráfica Q-Q esta gráfica se vería de la siguiente manera.



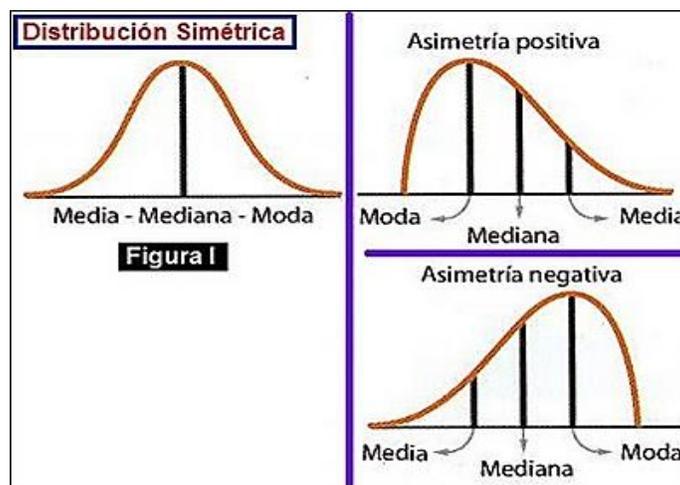
Nótese que, en esta gráfica, casi todos los puntos caen sobre la recta que representa los valores acumulados de una distribución normal teórica estandarizada.

Medidas de asimetría y curtosis de una distribución

Mediante la exploración gráfica, es posible visualizar si las distribuciones se apartan de la normal. El histograma y el polígono de frecuencia son los gráficos más sencillos para ver los sesgos de la distribución, y como vimos anteriormente, existen las gráficas P-P y Q-Q que cuantifican puntualmente la asimetría. Otras medidas de asimetría se

calculan directamente de la relación entre media, mediana y modo. Siendo la distribución normal simétrica (y por ello tanto el modo, la mediana y la media coinciden en la parte central de la distribución), la resta de las medidas mencionadas será cero. De esta propiedad se derivan dos coeficientes de asimetría conocidos como: a) Coeficiente de Asimetría de Pearson, y b) Coeficiente de Asimetría de Fisher.

En ambos casos, cuando el coeficiente de asimetría sea cero, indicará que el ajuste de la distribución observada es perfecto. En cambio, si el coeficiente de asimetría es mayor que cero ($Ca > 0$), indica que la distribución muestra un sesgo positivo. Al contrario, si el coeficiente de asimetría es menor que cero ($Ca < 0$), indica que la distribución muestra un sesgo negativo.



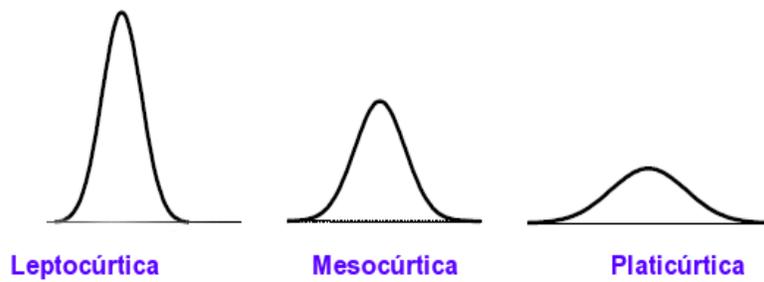
Otros coeficientes de asimetría toman en consideración la distribución de casos en cada uno de los cuartiles. En una distribución normal, la distribución de casos a lo largo de todos los valores de la variable será proporcional en cada uno de los cuartiles. Sobre esta base se calcula el coeficiente de asimetría de Bowley

$$Ab = \frac{(Q3 + Q1 - 2Me)}{(Q3 + Q1)}$$

En ambas expresiones Q son los cuartiles, y Me es la mediana de la distribución. En este caso también se cumple que si la distribución se ajusta a la normal, el valor del coeficiente es cero. Luego, será mayor que cero en el caso en que la distribución muestre un sesgo positivo, y menor que cero en el caso en que la distribución tenga asimetría negativa.

Otra medida que acompaña a las de simetría es la curtosis, y se refiere a la

acumulación de casos en torno a la media en una distribución. También toma como modelo a la distribución normal, cuyo valor de referencia en este estadístico es cero. La distribución normal por definición es mesocúrtica, por lo tanto el coeficiente de curtosis será cero para toda distribución que se empareje a la normal. Si la distribución empírica bajo análisis acumula la mayoría de los casos próximos a la media, su forma aparece como más puntiaguda que la normal. En tal caso se dice que la distribución es leptocúrtica y su coeficiente de curtosis será mayor que cero. Al contrario, si la distribución observada muestra que los casos se reparten más dispersos en torno a la media, se trata de una distribución platocúrtica y su coeficiente de curtosis será menor que cero. En el siguiente gráfico se muestran los distintos tipos de distribuciones.



Los coeficientes de asimetría y curtosis, nos indican el grado en que una distribución empírica de datos se aparta de una distribución normal. Cuando existen variaciones menores, evidenciada por estos coeficientes y los gráficos correspondientes, no existen restricciones para asimilar la distribución empírica y la normal estándar bajo el modelo matemático de esta última. Pero si las desviaciones aparecen marcadas es necesario recurrir a diferentes pruebas de normalidad estadística que no solo indican cuánto se aparta una distribución empírica de tal modelo, sino que tipo de correcciones deben hacerse. Estas pruebas son conocidas también como Bondad de Ajuste.

Puesto que el tema sobrepasa el desarrollo en este trabajo, solo mencionaremos algunas de tales pruebas. Todas se basan en el mismo principio: el valor de cada prueba cuantifica la probabilidad de que la distribución empírica observada provenga de una distribución normal. Si la probabilidad es alta, el modelo normal puede aplicarse, pero si tal probabilidad es baja, se descarta esa posibilidad y la distribución debe ser tratada bajo otro modelo matemático distinto que la normal estándar.

Pruebas de bondad de ajuste

Test de Shapiro-Wilk	Se usa para contrastar la normalidad de un conjunto de datos. Fue publicado por Samuel Shapiro y Martin Wilk. Se considera uno de los test más potentes para el contraste de normalidad, sobre todo para muestras pequeñas ($n < 50$).
Prueba de Kolmogórov-Smirnov	Es una prueba no paramétrica que determina la bondad de ajuste de dos distribuciones de probabilidad entre sí. Es una alternativa para verificar la normalidad de una distribución.
Prueba de Anderson-Darling	La prueba asume que no existen parámetros a estimar en la distribución que se está probando, en cuyo caso la prueba y su conjunto de valores críticos siguen una distribución libre. Cuando se aplica para probar si una distribución normal describe adecuadamente un conjunto de datos, es una de las herramientas estadísticas más potentes para la detección de la mayoría de las desviaciones de la normalidad.
Test de Jarque-Bera	En estadística, el test de Jarque-Bera es una prueba de bondad de ajuste para comprobar si una muestra de datos tiene la asimetría y la curtosis de una distribución normal. El test recibe el nombre de Carlos Jarque y Anil K. Bera.

Capítulo 6

Estimación de parámetros

Estudiar una población completa suele ser muy costoso, en tiempo, recursos materiales y humanos. Por ello, mediante un muestreo probabilístico de la población de interés, es posible estimar cualquier parámetro que estemos interesados conocer.

Cuando calculamos un índice en una población, lo denominamos parámetro. Cuando solo contamos con una muestra aleatoria de la población, y calculamos el mismo índice nos referimos a él como un estadístico. Si hemos asegurado un muestreo aleatorio de la población, es posible utilizar el índice calculado para estimar, con cierta confiabilidad, el parámetro poblacional. Como se desprende de lo dicho, la estadística nos sirve en este caso para conocer aspectos fundamentales de la población a un costo mucho menor. El proceso mediante el cual estimamos un parámetro a partir de una muestra se denomina inferencia estadística. Es por ello que para distinguir un parámetro de un estadístico se reserva el uso de las letras del alfabeto griego para referirnos a parámetros y del alfabeto latino para los estadísticos.

Estimación de una media poblacional

Suponemos que una variable cualquiera en la población tiene una distribución, que puede ser conocida o no. Mediante el teorema del límite central sabemos que sucesivas muestras del mismo tamaño extraídas de esa población, se aproxima a la distribución normal. Es decir, si procedemos a extraer muestras aleatorias de igual tamaño de la población (procedimiento que teóricamente puede repetirse infinitas veces), encontraremos que la distribución del estadístico en esa distribución sigue siempre un modelo normal.

Supongamos que deseamos conocer el promedio de años de estudios alcanzados en la población comprendida entre 30 y 40 años de edad residente en Córdoba Capital. El parámetro en el que estamos interesados es un promedio al que vamos a denominar μ (letra griega que se pronuncia mu), por tratarse de la población. Si pudiéramos extraer sucesivas muestras de esa población y calculamos el promedio de años de estudio en ella, encontraríamos que el estadístico sigue una distribución normal. Es decir, la mayoría de las muestras arrojarán un promedio o media, que se agruparán en torno a un valor particular. Lo que demuestra el teorema del límite central es que el valor del estadístico, en este caso la media, estará próximo al verdadero valor de μ , siempre que la muestra sea aleatoria.

Nótese que la media estará próxima pero no será exactamente igual a μ . Es decir, la

estimación de una media poblacional mediante una media muestral estará sujeta a error. Lo interesante del teorema del límite central es que especifica la distribución de ese error, llamado error estándar de la media, y es posible calcularlo mediante la siguiente ecuación:

$$\sigma_M = \frac{\sigma}{\sqrt{N}}$$

Conociendo las propiedades del teorema del límite central y la distribución de muestreo del estadístico, es posible estimar el parámetro estableciendo un límite de confianza, basado en las propiedades de la distribución normal estandarizada. Este proceso es el que se denomina estimación de parámetros.

Continuando con el ejemplo anterior, supongamos que se realiza un estudio tomando una muestra aleatoria de la población de 600 personas. Los resultados muestran que el promedio de años de estudio es 15,2 con una desviación estándar de 9.

Dado que se cuenta con una muestra representativa de la población de interés, se decide estimar el parámetro poblacional años promedio de estudios en la población adulta de 30 – 40 años, de la Ciudad de Córdoba Capital.

Datos:

Promedio de años de estudio: 15,2

Desvío Estándar: 9

Tamaño de la muestra: n=600

Para establecer los límites de confianza de la estimación, recurrimos a las áreas bajo la curva normal estandarizada. Sabemos que estas áreas se expresan en puntuación Z, dado que la distribución tiene como propiedad distintiva que su media es cero y su desviación estándar es igual a uno. En primer lugar necesitamos establecer el porcentaje de confiabilidad para la estimación. Supondremos para este estudio que queremos una confianza elevada para la estimación, esto es, del 99%. El paso siguiente es traducir ese porcentaje a valor Z. Por las propiedades de la normal estándar sabemos que un valor $Z = \pm 2,58$ cubre el 99% del área bajo la curva, de modo que ese dato es el que necesitamos para ajustar la estimación.

Puesto que solo estamos utilizando una sola muestra para estimar el parámetro, es necesario conocer el error estándar de la media para el intervalo de confianza. En la ecuación original, el error estándar se calcula a partir de la desviación estándar de la población que es σ (sigma). Como este es un parámetro y por tanto resulta desconocido, reemplazamos en la fórmula original estos términos por la desviación estándar muestral y a N por n-1 (tamaño muestral menos uno). Así el error

estándar para la media de este ejemplo será:

$$EE_M = \frac{DE}{\sqrt{n-1}} = \frac{9}{\sqrt{199}} = 0,637$$

Ahora debemos calcular el Intervalo de Confianza para la media. El intervalo establece dos límites, uno superior y otro inferior, dentro de los cuales se supone que se encuentra el parámetro poblacional. La ecuación para calcular ese intervalo es la siguiente:

$$LI = M \pm (Z \times EE_M)$$

En nuestro ejemplo estos límites serán

Superior: $15,2 + (2,58 \times 0,367) = 16,14$

Inferior: $15,2 - (2,58 \times 0,367) = 14,25$

El resultado se expresa de la siguiente manera: $14,25 \leq \mu \leq 16,14$. Así, es posible afirmar que la cantidad de años de estudio promedio de la población de Córdoba Capital, entre 30 y 40 años, se encuentra comprendida entre los 14 y los 16 años aproximadamente (con una confianza del 99%).

Estimación de una proporción poblacional

De la misma manera en que se puede estimar una media paramétrica, es también posible estimar otros parámetros en la población, como por ejemplo una proporción. La distribución de muestreo de una proporción sigue los mismos principios que la distribución para la media, es decir, que se pueden explicar mediante el teorema del límite central, y por tanto la distribución de referencia será la normal estandarizada.

Veamos mediante un ejemplo cómo se estima una proporción paramétrica. Supongamos que se desea estimar la tasa de promoción entre primaria y secundaria en la ciudad de Córdoba Capital. Para ello se toma una muestra de 500 alumnos de diez escuelas de la ciudad. Se encuentra que en la muestra, 439 alumnos son efectivamente promovidos. La tasa resultante es $439/500=0.87$. Con estos datos se estima que la tasa de promoción en la población escolar en la ciudad de Córdoba Capital es del 87%.

Datos

Proporción de alumnos promovidos: $p=0.87$ (87%)

Tamaño de la muestra: 500

Como en el caso anterior, es necesario fijar un nivel de precisión para la estimación. En este caso el nivel de precisión será del 95%. Aquí nuevamente consultamos la distribución normal estandarizada y obtenemos el valor Z que corresponde al 95% del área bajo la curva que es igual a $\pm 1,96$.

Con los datos obtenidos debemos calcular el error estándar de la proporción, de la misma manera en que calculamos el error estándar de la media en el ejercicio anterior. La ecuación de cálculo es la siguiente:

$$EE_p = \sqrt{\frac{p \times (1 - p)}{n}}$$

Reemplazando los términos en la fórmula tenemos que el error estándar de la Proporción es:

$$EE_p = \sqrt{\frac{0.87 \times (1 - 0.87)}{500}} = 0,015$$

Para calcular el Intervalo de Confianza para la proporción, utilizamos la misma ecuación que la empleada para el intervalo de confianza de la media, reemplazando en ella los términos correspondientes. Así tenemos que en este caso la ecuación se define como:

$$LI = p \pm (Z \times EE_p)$$

Así, los intervalos de confianza buscados son:

Superior: $0,87 + (1,96 \times 0,015) = 0,8994$

Inferior: $0,87 - (1,96 \times 0,015) = 0,855$

Obtenemos como resultado lo siguiente; $85,5\% \leq p \leq 89,94\%$. Ello se interpreta de la siguiente manera: la tasa de promoción entre primaria y secundaria en la ciudad de Córdoba Capital se encuentra comprendida entre el 90% y el 85,5% aproximadamente.

Comparación de proporciones

Con frecuencia ocurre que se toman estadísticas continuas y es posible obtener datos poblacionales históricos que nos sirven de referencia para indagar en qué medida la situación reflejada en ellos pudiera haber cambiado. Para explicar sobre esta situación sirvámonos del siguiente ejemplo: De acuerdo a datos tomados en el año 2001, se constató que el 67% de los puestos gerenciales de las principales empresas del país estaba ocupado por varones. En el año 2010, se realiza un nuevo estudio tomando una muestra aleatoria de 249 gerentes de las principales empresas del país, de los cuales la proporción de varones es de 64,66%. Un investigador social se pregunta si existen diferencias significativas entre las proporciones encontradas en estos dos periodos respecto al género y el cargo dentro de las empresas. Para resolver estadísticamente este problema, primero debemos reducir el problema a proporciones, tal como se definen para el estudio realizado en 2010, en tal caso tenemos que:

Proporción de varones en la muestra: $p=0.6466$ (64,66%)

Proporción de mujeres en la muestra: $q= 1-p= 1- 0.6466= 0.3534$ (35,34%)

Se trata de ver ahora si la proporción encontrada en 2001 es diferente a la encontrada en 2010. La prueba de proporciones consiste en comparar un valor Z empírico, con un valor Z teórico definido a partir de la probabilidad asignada a la diferencia. El valor de Z empírico se obtiene aplicando la siguiente fórmula:

$$Z_{empirico} = \frac{p_{muestra} - p_{población}}{\sqrt{\frac{p_{población} - q_{población}}{n_{muestra}}}}$$

Reemplazando los valores del problema en la fórmula tenemos que:

$$Z_{empirico} = \frac{0,6466 - 0,67}{\sqrt{\frac{0,67 - 0,33}{249}}} = -0,65$$

Para el ejemplo que estamos desarrollando, el valor de Z empírico es -0.65. El valor Z

teórico definido a partir de la probabilidad asignada a la diferencia, es el equivalente al intervalo de confianza para la estimación del parámetro. En este caso el valor Z teórico se asigna con anterioridad al desarrollo de la investigación, pero a los fines de este ejemplo lo asignamos en este punto. Digamos que deseamos tener una seguridad del 99% para determinar si existe una diferencia entre las proporciones. Por lo ya visto, sabemos que el valor Z teórico correspondiente a ese valor de área bajo la curva de distribución normal es igual a $Z \pm 2,58$. Por lo tanto para aseverar que la diferencia entre las proporciones es significativa, el valor de Z empírico, deberá ser mayor que 2,58 o menor que -2,58. Solo aquellos valores que se ubiquen en los extremos definidos por el valor Z teórico dado, serán considerados suficientes como para determinar que la diferencia es estadísticamente significativa. En este ejemplo tenemos que $Z_{empirico} - 0.65 > Z_{teórico} - 2,58$, por lo cual no es posible aceptar que existe una diferencia de proporciones y que la diferencia observada se debe a la aleatoriedad de los datos.

Comparación de proporciones muestrales

En el caso anterior se contaba con una proporción encontrada en una población y se la comparó con la proporción encontrada en una muestra. Cuando solo se cuenta con proporciones extraídas de muestras, igualmente se puede determinar si una diferencia de proporciones resulta significativa. En tal caso se utiliza una distribución teórica muestral de diferencia de pares de proporciones de muestreo, pertenecientes a una misma población. Esta distribución, tal como hemos visto hasta aquí, también sigue un modelo normal, y en este caso, un supuesto central para esta distribución teórica es que la media de las diferencias de proporciones es igual a cero. Bajo el mencionado supuesto, el error estándar de la distribución muestral de diferencias de proporciones será igual a:

$$\sigma_{dif} = \sqrt{\frac{p_1 \times q_1}{n_1} + \frac{p_2 \times q_2}{n_2}}$$

Veamos un ejemplo de la manera en que procederíamos al comparar dos proporciones muestrales. Supongamos que mediante un muestreo aleatorio de la población, se cuenta con dos proporciones referidas a la tasa de desempleo en las ciudades de Paraná y Concordia, tomando como referencia que ambas forman parte de la población de la Provincia de Entre Ríos. La muestra 1 ($n_1=500$) tomada de Paraná muestra que la tasa de desempleo es de 15% (p_1). La muestra 2 tomada en Concordia ($n_2=280$) muestra la tasa de desempleo es del 25% (p_2). Nos interesa determinar si las

dos proporciones encontradas difieren significativamente, o si la diferencia observada es debida al azar. Para ello debemos establecer un intervalo de confianza donde situar la diferencia encontrada. Como ya sabemos que la diferencia de pares de proporciones de muestreo pertenecientes a una misma población sigue una distribución normal, el valor del intervalo lo establecemos utilizando los valores de área bajo la curva normal estandarizada. Para este ejemplo, vamos a establecer un intervalo de confianza del 95%, por lo tanto el valor crítico de Z será ± 1.96 . Ahora procedemos a calcular el error estándar de la diferencias de proporciones:

$$\sigma_{dif} = \sqrt{\frac{0,15 \times 0,85}{500} + \frac{0,25 \times 0,75}{280}} = 0,0179$$

Necesitamos ahora calcular el valor Z empírico, con el cual compara el valor Z crítico establecido en el intervalo de confianza elegido. Para esto debemos resolver la siguiente ecuación:

$$Z_{empírico} = \frac{p_1 - p_2}{\sigma_{dif}}$$

Con los datos que se tienen, el valor Z empírico es igual a

$$Z_{empírico} = \frac{0,15 - 0,25}{0,0179} = -5,58$$

Dado que $-5,58 < -1.96$ se puede concluir que la diferencia de proporciones es significativamente diferente. Es decir, no son producto de las variaciones debidas al azar o error de muestreo.

En este ejemplo debemos tener en cuenta dos cuestiones importantes, que por la simplicidad de la presentación hemos dejado de lado. Primero; lo esbozado en los párrafos precedentes tiene sentido cuando se comparan proporciones extraídas de la misma población, dado que este es un supuesto central para aplicar la distribución de muestreo de diferencia de proporciones. En este caso, la población de referencia de donde se extraen las muestras es la provincia de Entre Ríos. Segundo; es necesario definir con anterioridad el sentido de la comparación, dado que en la segunda ecuación, debemos calcular la diferencia de proporciones según lo que hayamos

establecido previamente como p_1 y p_2 . Generalmente, se parte de algún supuesto de base para establecer qué proporción es cada una, de lo contrario perdería todo sentido la comparación.

Estimación de parámetros: conceptos teóricos

Estimación puntual: en la estimación puntual se da como valor estimado de μ el valor de la media muestral. De acuerdo al teorema del límite central el valor de la media muestral será igual o estará muy próximo a μ . La diferencia entre la media muestral y μ , corresponde al error propio del muestreo, pero tal diferencia será pequeña. Entonces, en una estimación puntual se dirá que el valor de μ estará próximo a la media muestral, o bien que la diferencia entre μ y la media muestral será pequeña.

Estimación por intervalos: en la estimación por intervalos, se establecen límites de confianza entre los cuales es más probable que se encuentre la media poblacional. En otras palabras, los límites determinan la probabilidad dentro de los cuales μ esté contenido. Esta probabilidad se llama probabilidad de confianza y el intervalo marcado por los límites se llama intervalo de confianza.

Dado que la distribución de muestreo de las medias es normal, se sostiene que si la media poblacional es μ , existe una probabilidad del 95% de que la media de la muestra aleatoria extraída de esa población tome un valor $\mu - (1.96 * \sigma_M)$ y $\mu + (1.96 * \sigma_M)$. Dicho de otra manera; el 95% de las muestras aleatorias extraídas de la población, tendrán una media comprendida entre $\mu - (1.96 * \sigma_M)$ y $\mu + (1.96 * \sigma_M)$.

Si se cuenta con una muestra aleatoria de la población, el intervalo de confianza para μ , se establece de la siguiente manera: $M - (1.96 * EEM)$ y $M + (1.96 * EEM)$.

Ello implica que:

- a) El intervalo de confianza está centrado en la media observada en una muestra aleatoria extraída de la población.
- b) Existe un 5% de probabilidades de que ese intervalo no contenga al parámetro μ . Dado que la probabilidad de ocurrencia es baja (5 de cada 100), estaremos en condiciones de afirmar que en la mayoría de los casos el intervalo de confianza contiene el valor de μ .
- c) En raras ocasiones se conocerá el error estándar de la media o σ_M , por lo que también deberá ser estimado a partir de la muestra. Para ello se utiliza la desviación estándar de la muestra, tal como mostró en los ejemplos.

Capítulo 7

Prueba de hipótesis sobre la media paramétrica

La investigación científica se inscribe fuertemente en lo que se denomina prueba de hipótesis. Es un tema que abarca numerosos capítulos de los manuales de metodología de la investigación, por lo que en este apartado solo haremos referencia a la manera de plantear una hipótesis de investigación que pueda ser puesta a prueba con un modelo estadístico. Es este otro aspecto de lo que denominamos estadística inferencial, que se refiere a la manera en que los datos empíricos pueden ser analizados con diferentes modelos matemáticos. De los muchos modelos estadísticos utilizados para la prueba de hipótesis nos centraremos en dos: el modelo normal y el modelo chi cuadrado. Ambos tienen la propiedad de que son excelentes ejemplos para generalizar a cualquier otro modelo estadístico.

Comencemos con lo básico que es comprender la estructura de una hipótesis. Una hipótesis se parece a lo que conocemos como una conjetura o suposición, pero tiene la propiedad de explicitar en qué medida es posible demostrar con cierto grado de rigor su veracidad o falsedad. Las hipótesis estadísticas son un caso especial de lo que mencionamos anteriormente ya que siempre es posible expresarlas de manera tal que, mediante una colección de datos, podamos ponerlas a prueba. Veamos un ejemplo: un grupo de investigadoras verifica que el promedio en matemáticas de las pruebas ONE 2013 es más alto que el obtenido en las pruebas ONE 2010 para los alumnos de Córdoba Capital. En base a esta información pueden suponer que el promedio del operativo de evaluación Aprender 2016 será también más alto. Esta suposición puede tomar la forma de una hipótesis estadística si la redactamos de la siguiente manera: el promedio de los alumnos de Córdoba Capital en el operativo nacional de evaluación Aprender 2016, será mayor que el obtenido en el ONE 2013. También podemos expresarlo estadísticamente:

$$H_1: \mu_{\text{aprender-2016}} > \mu_{\text{ONE-2013}}$$

El planteo estadístico de la hipótesis de investigación se lee del siguiente modo: la media poblacional en el operativo de evaluación Aprender 2016, será mayor a la media poblacional obtenida en ONE 2013. Esta hipótesis es válida para la población de escolares de Córdoba Capital y para el área de matemáticas.

Las investigadoras también podrían haber hecho otra suposición si los datos fueran diferentes. Tomemos el caso de que los promedio obtenidos por los escolares de Córdoba Capital para las evaluaciones de matemáticas de ONE 2010 sean un poco más

bajas que las obtenidas en el ONE 2013. Las investigadoras ahora podrían suponer que habrá diferencias más grandes entre el ONE 2013 y el operativo aprender 2016. El problema que enfrenan es que no pueden decidir el sentido de la diferencia, tal como se planteaba en la hipótesis anterior. En tal caso la hipótesis queda redactada de la siguiente manera: el promedio de los alumnos de Córdoba Capital en el operativo nacional de evaluación Aprender 2016, será diferente que el obtenido en el ONE 2013. También podemos expresarlo estadísticamente:

$$H_1: \mu_{\text{aprender-2016}} \neq \mu_{\text{ONE-2013}}$$

El planteo estadístico de la hipótesis de investigación se lee de la siguiente manera: la media poblacional en el operativo de evaluación Aprender 2016, será diferente a la media poblacional obtenida en ONE 2013. Nuevamente esta hipótesis es válida para la población de escolares de Córdoba Capital y para el área de matemáticas.

Podemos establecer cualquier hipótesis respecto de las diferencias en los promedios o proporciones. Incluso podemos establecer hipótesis sobre la magnitud de las correlaciones y asociaciones entre variables. En estadística, las hipótesis plantadas siempre descansan en un modelo matemático que permite ponderar las diferencias o magnitudes y así decidir sobre la validez de la hipótesis.

Los datos y el modelo estadístico

Antes de continuar con el planteo de hipótesis, debemos adentrarnos al tema de los modelos estadísticos. Intentaremos hacerlo apelando a la mayor simplicidad matemática sin perder la rigurosidad.

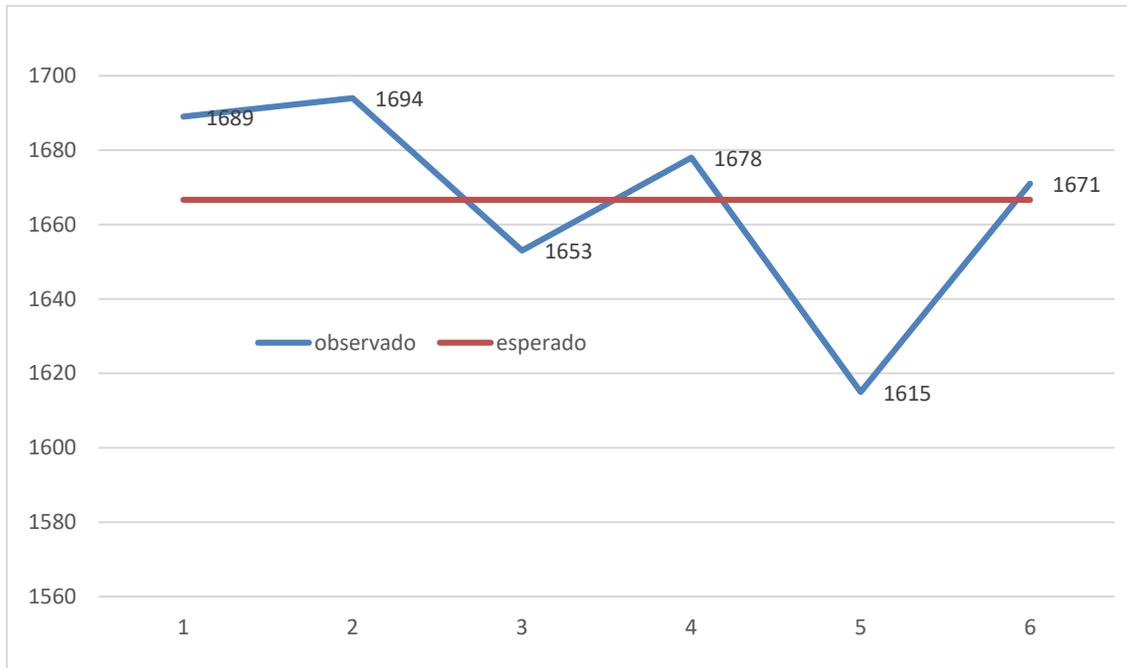
Comencemos con una distribución de datos teórica que proviene del lanzamiento de un dado. El lanzamiento del dado puede considerarse como un espacio de posibilidades perfectamente conocido, pues sabemos que el resultado variará entre 1 y 6. Suponiendo que el dado no tiene ninguna falla, si lanzamos una serie infinita de veces el dado el valor de probabilidad para cada cara del dado es igual a $1/6$ o bien, $p(x)=0,16$, siendo x cualquier valor entre 1 y 6. En la siguiente tabla vemos el valor de x que varía entre 1 y 6, luego la probabilidad de que se obtenga un valor dado en un lanzamiento del dado para un espacio muestra igual a: $E\{1,2,3,4,5,6\}$.

Sabemos de antemano que esa posibilidad es la probabilidad de x , que es igual a $p(x)=1/6$. La columna $f(x)$ contiene la frecuencia de una simulación de diez mil tiradas de un dado bajo el espacio muestral $E\{1,2,3,4,5,6\}$, con una probabilidad de $1/6$ para cada evento. La columna $f(x)p(x)$ es la frecuencia esperada en una distribución de 10.000 tiradas con la probabilidad exacta de que cada resultado ocurra $1/6$ de las veces.

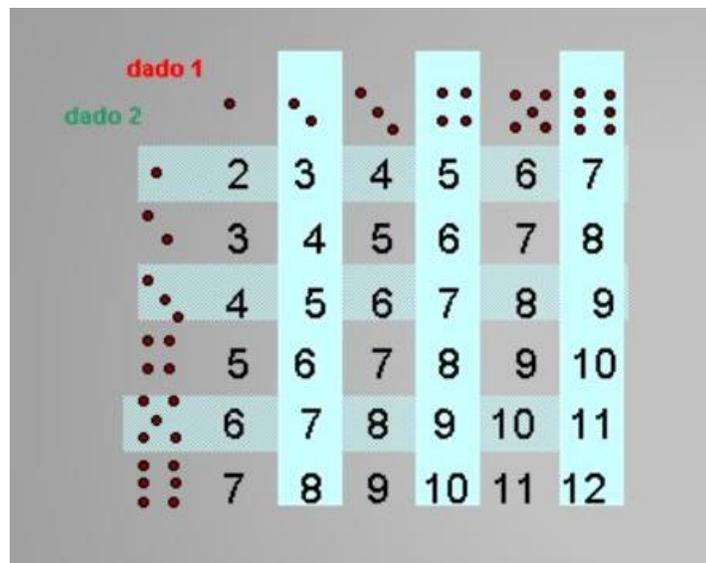
Valor	$p(x)$	$f(x)p(x)$	$f(x)$
1	1/6	1666.66	1639
2	1/6	1666.66	1584
3	1/6	1666.66	1633
4	1/6	1666.66	1638
5	1/6	1666.66	1577
6	1/6	1666.66	1621

Entonces, si pudiéramos repetir el evento una cantidad infinita de veces tendríamos una distribución de datos perfectamente uniforme, pues sabemos que la ley de los grandes números determina que todo proceso aleatorio (en este caso la tirada de un dado), se empareja a su probabilidad de ocurrencia (en este caso, que cada valor aparezca 1/6 de las veces). Si tiramos el dado unas diez mil veces esperamos que cada valor de $E\{1,2,3,4,5,6\}$, se repita unas 1666 veces.

Sin embargo, las distribuciones empíricas (en este caso los resultados de la simulación), se aproximan a los valores esperados, pero no se igualan perfectamente. En la siguiente gráfica representamos los valores de $f(x)$ y $f(x)p(x)$; $f(x)$ representa las frecuencias observadas en la modelización para una secuencia de 10.000 tiradas de un dado, $f(x)p(x)$ es el modelo teórico de las frecuencias esperadas para una secuencia de 10.000 tiradas de un dado. Vemos que existen diferencias entre el valor esperado y el observado, y a esto se lo llama fluctuación aleatoria. La fluctuación aleatoria ocurre porque la simulación solo toma un número finito de observaciones, en este ejemplo fueron 10.000, y el valor esperado se cumple cuando el valor es infinito. Es de notarse que, a pesar de las fluctuaciones en torno al valor esperado, éstas no son tan diferentes; todas están próximas a 1.666. Sobre los datos que hemos presentado podemos afirmar que los valores observados pertenecen a una distribución uniforme del espacio muestral $E\{1,2,3,4,5,6\}$ con una probabilidad 1/6 para cada evento.



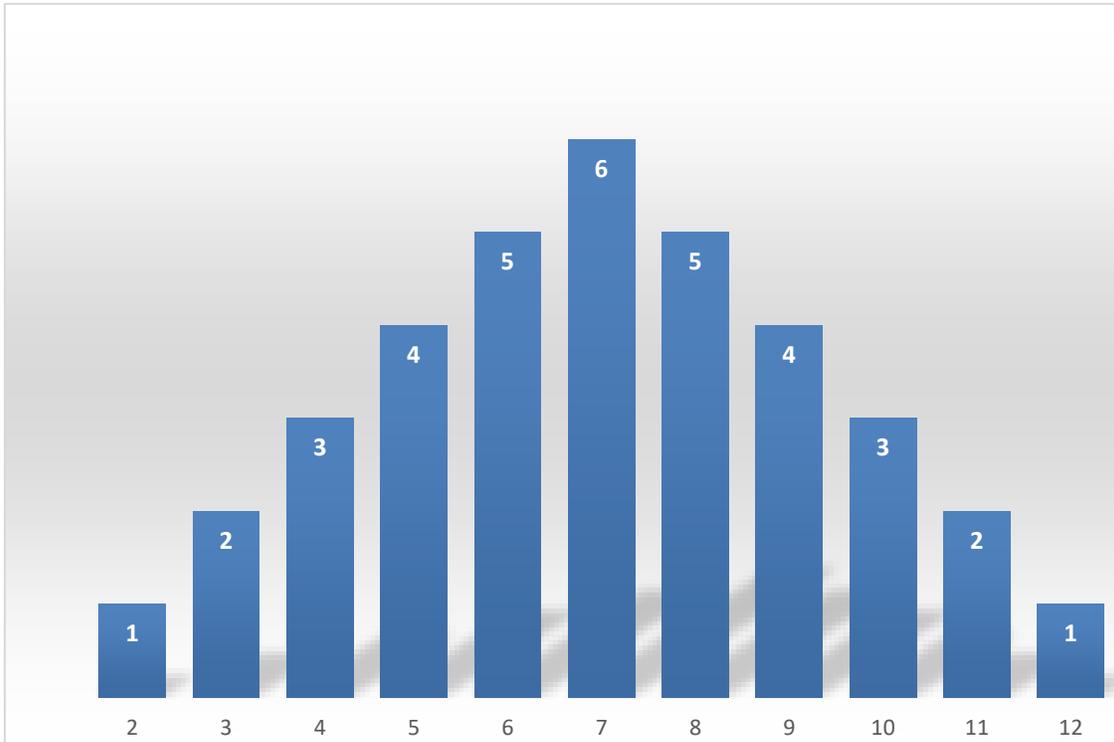
Ahora, repetiremos todo el proceso anterior, pero teniendo en cuenta el resultado del lanzamiento de dos dados. El espacio muestral sobre el que ahora trabajamos resulta de la suma de los valores de las dos caras del dado, por lo tanto, tenemos una variación de 2 a 12. Es decir, $E\{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$. Ahora, cada uno de los valores de E no tiene la misma posibilidad de ocurrir; solo hay una forma en de obtener un 2 en la suma de las caras del dado, y hay seis formas de obtener un 7 en la suma de las caras del dado. Para ilustrar la combinatoria de resultados posibles véase la siguiente figura.



En la primera fila están las seis caras del dado 1, en la primera columna las seis caras del dado 2. Lo que buscamos es el resultado de la suma de los valores de las caras de ambos dados, por lo cual existen 36 resultados posibles. Si observamos la cuadrícula, vemos que solo hay una forma de obtener un 2, que es cuando se obtiene 1 en ambas caras de los dados. Luego vemos que hay dos formas de obtener un 3 que es cuando obtenemos 1 en la cara del primer dado y obtenemos 2 en la cara del segundo dado, y a la inversa. Siguiendo este razonamiento vemos que hay tres formas de obtener un 4 y así para los demás resultados. Por lo tanto, es factible construir una tabla de frecuencias con las probabilidades de ocurrencia de cada resultado de la suma de las caras del dado, que denominaremos x .

Valor de x	$f(x)$	$p(x)$
2	1	1/36
3	2	2/36
4	3	3/36
5	4	4/36
6	5	5/36
7	6	6/36
8	5	5/36
9	4	4/36
10	3	3/36
11	2	2/36
12	1	1/36

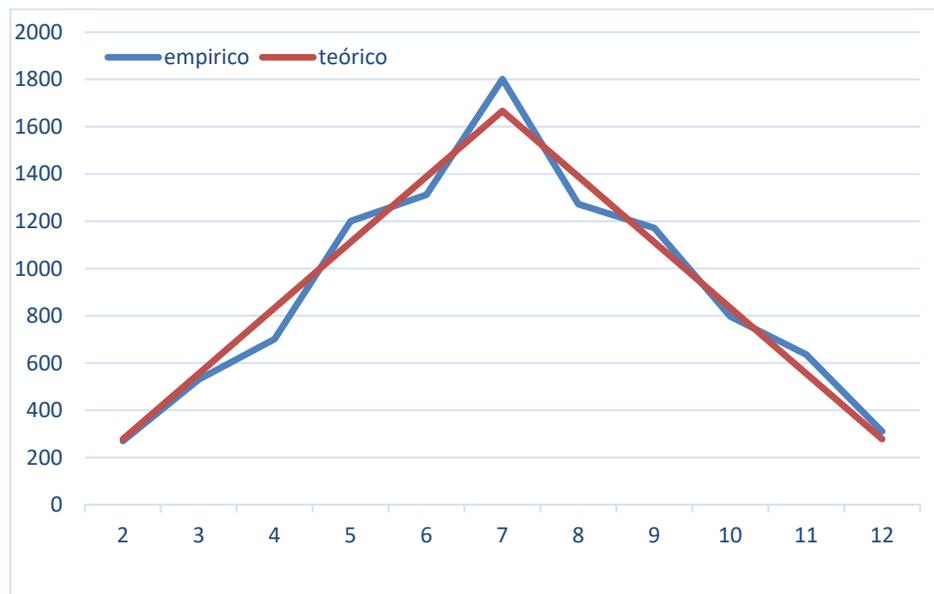
La gráfica resultante de esas frecuencias se muestra a continuación y en ella podemos ver qué valores de la suma de las dos caras de un dado tienen mayores chances de ocurrir. Dicho de otro modo, la gráfica nos muestra que hay seis combinaciones posibles para obtener un 7, hay cinco combinaciones posibles para un 6 y 8, y así para los otros resultados tenemos que existen menos combinaciones posibles.



Anteriormente, conociendo la probabilidad de ocurrencia del valor de la tirada de un dado obtuvimos una distribución uniforme que empleamos para simular los resultados de una tirada de 10.000 veces un dado. Podemos repetir el mismo proceso de simulación para la tirada de dos dados, basándonos en las probabilidades teóricas de obtener cualquier valor de $E\{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$. En la tabla, esas probabilidades se representan por $p(x)$ como vimos anteriormente, ese el modelo probabilístico esperado de la distribución de frecuencias de la suma de las dos caras de un dado en una tirada de diez mil veces. Luego $f(x)p(x)$ son las frecuencias efectivamente observadas en la simulación de la suma de las dos caras de un dado en una tirada de diez mil veces.

	$p(x)$	$f(x)p(x)$
2	277,7	270
3	555,5	530
4	833,3	801
5	1111,1	1200
6	1388,8	1312
7	1666,6	1801
8	1388,8	1272
9	1111,1	1171
10	833,3	797
11	555,5	536
12	277,7	310

Al graficar los valores de la tabla vemos que el modelo teórico es una distribución perfectamente simétrica con su valor máximo en el resultado $x=7$. La distribución empírica se aproxima bien a ese modelo teórico a pesar de las fluctuaciones aleatorias.



Lo que hemos visto hasta aquí fueron dos simulaciones de distribuciones de datos basadas en un modelo probabilístico conocido. Las distribuciones observadas se aproximaban mucho a los modelos utilizados, y esta diferencia entre el modelo teórico y los datos empíricos le llamamos fluctuación aleatoria. Otra manera de expresar lo dichos es afirmando que una distribución de datos empírica puede modelizarse mediante una distribución de probabilidades conocida. Cuando ello ocurre decimos

que la distribución de datos empírica se ajusta a un modelo teórico y que las fluctuaciones aleatorias son mínimas. Es así que las propiedades matemáticas del modelo teórico pueden utilizarse para describir el conjunto de datos empíricos. Existen muchos modelos estadísticos que sirven a ese propósito, es por ello que para cada conjunto de datos se estima una bondad de ajuste a un modelo dado.

La distribución normal como modelo estadístico

En los ejemplos anteriores establecimos previamente un modelo probabilístico y procedimos a simular un conjunto de datos tomando como referencia el modelo propuesto. En el apartado en que describimos la distribución normal, se mencionó que ésta podía ser utilizada como modelo estadístico y que es uno de los más usados en estadística.

Comencemos describiendo un conjunto de datos $E\{2, 4, 6\}$. Para este conjunto calculamos las medidas de tendencia central y las de dispersión, que son las siguientes:

promedio: 4

Varianza: 2.66

Desviación Estándar: 1.63

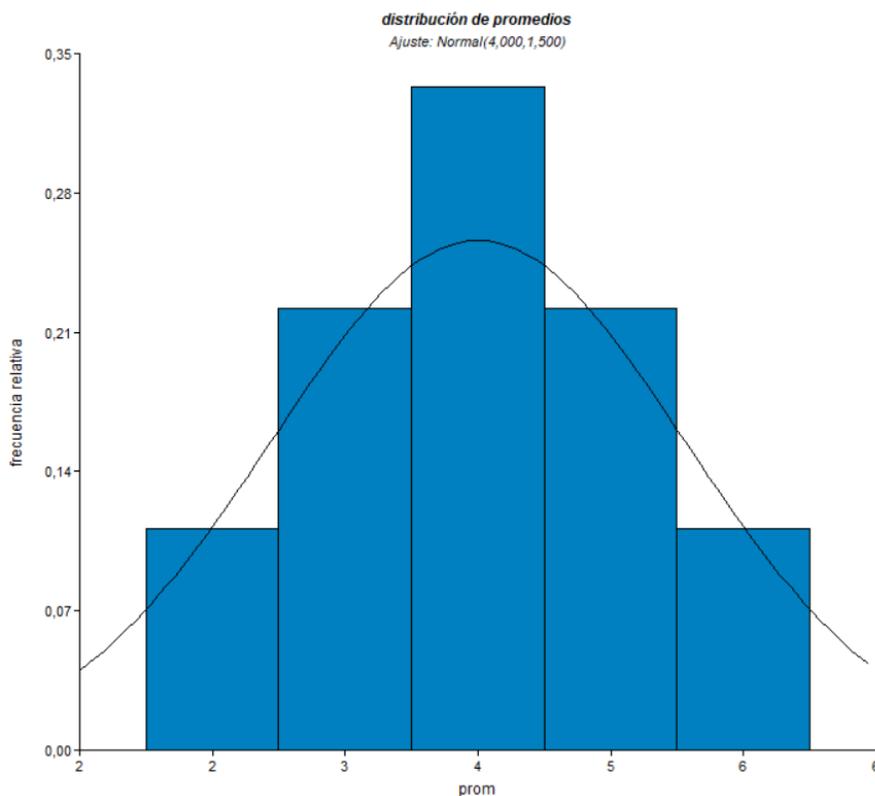
Ahora, vamos a calcular todos promedios del conjunto $E\{2, 4, 6\}$, pero tomado valores de a pares. Es decir, vamos a generar una distribución de promedios a partir del conjunto $E\{2, 4, 6\}$. Como en el ejemplo de los dados, es posible pasar los valores a una tabla cuyas filas y columnas sean los valores del conjunto $E\{2, 4, 6\}$, y las celdas contengan los promedios de los pares de valores.

Distribución de promedios	2	4	6
2	2	3	4
4	3	4	5
6	4	5	6

Esta información la pasamos a una tabla que contenga la frecuencia de con la que ocurre cada promedio y la probabilidad de obtener cada uno de ellos.

Promedio	f(x)	p(x)
2	1	1/9
3	2	2/9
4	3	3/9
5	2	2/9
6	1	1/9

Ahora podemos trazar un histograma de la distribución de promedios y comprobar empíricamente a qué tipo de distribución probabilística podemos asimilarla.



Vemos que la distribución de los promedios del conjunto $E\{2, 4, 6\}$ da como resultado un gráfico que se asemeja a una distribución normal. Es decir, la distribución de promedios del conjunto $E\{2, 4, 6\}$ tomado de a dos valores, produce una distribución normal. Esta distribución está centrada en el valor promedio 4, siendo menos frecuentes los promedios 2 y 5, y menos frecuente aún los promedio 2 y 6. Esto se cumple para cualquier distribución de promedios de una variable aleatoria y ha sido descrito como teorema del límite central. Este teorema tiene toda una formulación matemática, pero intentaremos explicarlo de manera sencilla. Nuestro conjunto de

datos original es $E\{2, 4, 6\}$ cuyo promedio es 4. Si procediéramos a seleccionar aleatoriamente dos valores cualquiera de ese conjunto vemos que es más probable escoger algunos valores que otros.

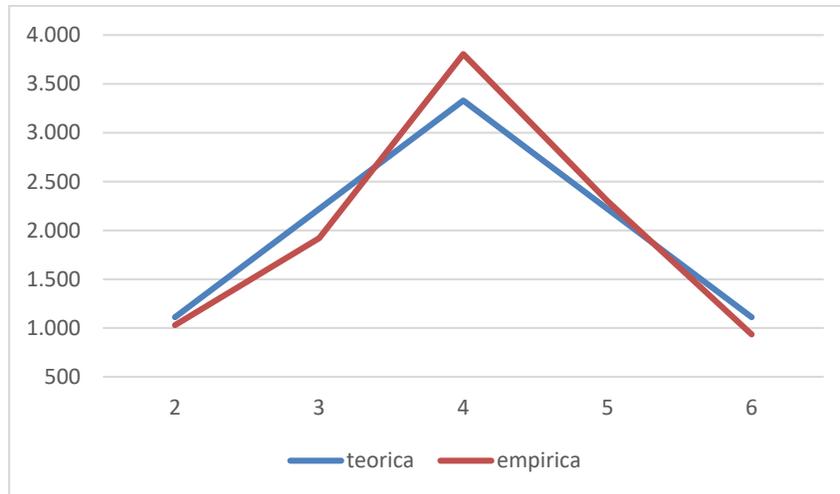
Para simplificar esto, diremos que $p(x)$ es la probabilidad de obtener un promedio dado y lo expresamos en valor porcentual. Luego tabulamos la distribución teórica dado los valores probabilísticos asignados de obtener cualquier promedio del conjunto $E\{2, 4, 6\}$ cuyos valores son tomados de a pares en una simulación de un muestro aleatorio de 10.000 ensayos. El resultado se muestra en la siguiente tabla:

Promedio	$p(x)$	$f(x)$ teórica	$f(x)$ empírica
2	11%	1.110	1030
3	22%	2.220	1920
4	33%	3.330	3805
5	22%	2.220	2301
6	11%	1.110	934

Según se lee en la tabla, la probabilidad de obtener un promedio de 2 en este experimento ficticio es solo del 11%, mientras que obtener un promedio de 4 tiene una chance del 33%. Teóricamente en un ensayo de 10.000 repeticiones, obtendríamos 1.100 muestras con promedio de 2 y unas 3.300 muestras con promedio de cuatro.

Existe una sola posibilidad en que podamos obtener un promedio de dos, y es que en el par de valores extraído del conjunto $E\{2, 4, 6\}$, ambos tengan valor 2. Pero la combinación de pares para obtener un promedio de 4 es mayor, dado que esto ocurre cuando la muestra de pares de valores de E se compone de $6+2$, $2+6$ y $4+4$. Entonces, el proceso de repetir el muestreo aleatorio de promediar pares de valores extraídos del conjunto E tiende a centrarse en 4, y los promedio 2 y 6 que solo pueden ocurrir por la combinación de dos valores, tienen a ser menos frecuentes y ocupar los extremos de la distribución. Es por ello que el modelo estadístico que mejor refleja el procedimiento es la distribución normal. La columna final tiene los datos de una simulación de 10.000 promedios basados en la combinación de dos valores del conjunto $E\{2, 4, 6\}$. Como vemos, se cumple que en el largo plazo los valores empíricos se aproximan a los valores empíricos con bastante fiabilidad.

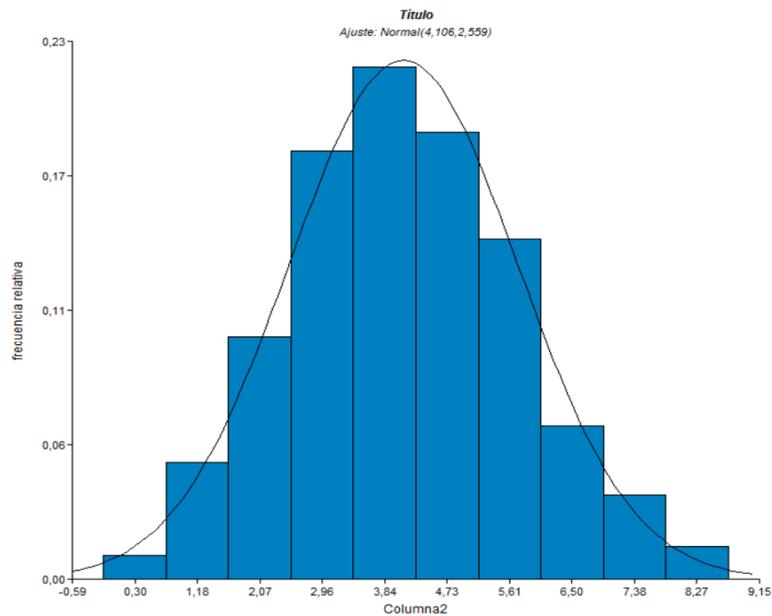
En el siguiente gráfico se muestran las curvas resultantes del modelo teórico de distribución de frecuencias, y el resultado empírico de una simulación de 10.000 promedio del conjunto $E\{2, 4, 6\}$ basados en $p(x)$. Nuevamente vemos la semejanza de las distribuciones y las fluctuaciones aleatorias entre el valor teórico y el empírico.



Sigamos con nuestro hipotético conjunto de datos $E\{2, 4, 6\}$, y supongamos que fue extraído por medio de un muestreo aleatorio de una población con los siguientes parámetros:

μ	σ
4,11	1,6

Sabiendo que el conjunto $E\{2, 4, 6\}$ tiene un promedio de 4, concluimos que se encuentra muy próximo al centro de la distribución y no hay razón para dudar de ello. Si observamos el gráfico de la población de origen del conjunto $E\{2, 4, 6\}$ comprobaremos de manera intuitiva lo dicho anteriormente.



Comparación de distintos promedios con el modelo estadístico de la distribución normal estándar

Vamos a tomar ahora tres conjuntos de datos e intentaremos decidir si pertenecen a una población con parámetros

μ	σ
4,11	1,6

Los conjuntos son los siguientes,

Conjunto A	Conjunto B	Conjunto C	Conjunto D
3,28	2,77	6,93	-3,5
6,7	8,05	7,05	-7,2
2,11	4,99	8,09	6,8
$\bar{x}=4,03$	$\bar{x}=5,27$	$\bar{x}=7,36$	$\bar{x}=-1,3$

Ordenando los conjuntos por el valor de su promedio, vemos que, de A a D, se alejan cada vez más del valor paramétrico. Ahora bien, si consideramos que en la distribución normal estandarizada los límites $Z=\pm 1,96$ cubren el 95% de la distribución, es posible suponer que cualquier valor de Z por encima o por debajo de $Z=\pm 1,96$ puede considerarse como de una ocurrencia poco probable, y si este es el caso, es lícito suponer que esa muestra: a) no pertenece a la distribución original, o b) es una muestra posible pero atípica para la distribución original.

Para considerar la posición de cada conjunto en la muestra original, vamos a transformar los valores promedios a puntuación Z y comprobar si se encuentran dentro del área comprendida por los valores $Z=\pm 1,96$.

Conjunto A	Conjunto B	Conjunto C	Conjunto D
$\bar{x}=4,03$	$\bar{x}=5,27$	$\bar{x}=7,36$	$\bar{x}=-1,3$
$Z=-0,05$	$Z=0,725$	$Z=2,03$	$Z=-3,38$
$p=0,961$	$p=0,469$	$p=0,043$	$p=0,001$

Como puede observarse, el conjunto A y B se encuentran dentro de los límites del valor $Z=\pm 1,96$. Por lo tanto no hay razón para dudar de que pertenecen al conjunto original de datos con parámetro $\mu=4,11$. En cambio, el conjunto C y D están fuera de los límites de $Z=\pm 1,96$. Recordemos una vez más que en una distribución normal estándar, el valor $Z=\pm 1,96$ cubre el 95% de la distribución, por lo tanto, los promedios de los

conjuntos C y D están por fuera de esa área. En otras palabras, los consideramos ahora como conjuntos atípicos para una población con $\mu=4,11$, y es lícito suponer que no pertenezcan a esa población.

Para afirmar lo dicho anteriormente hemos agregado a la tabla la probabilidad de ocurrencia de una muestra con los promedios de los conjuntos A, B, C y D, si consideramos que efectivamente todos ellos pertenecen a una población con $\mu =4,11$. Para una interpretación más directa de la probabilidad, transformamos los valores a porcentajes y tenemos que para el conjunto A la probabilidad de pertenecer a la población de origen es de 96,1%, para el conjunto B la probabilidad de pertenecer a la población de origen es 46,9%. Ahora la probabilidad del conjunto C de pertenecer a la población de origen es de 4,3% y la del conjunto D solo del 0,1%.

Supongamos ahora que tenemos que responder a la siguiente pregunta: ¿dado los valores del conjunto C, es posible afirmar que pertenecen a una población con $\mu=4,11$? Dado que la probabilidad de ocurrencia de una muestra de tales características es menor al 5%, es posible afirmar que ese conjunto de valores pertenece a otra población.

Prueba de hipótesis sobre media paramétrica.

Ahora vamos a utilizar todos los conceptos descriptos anteriormente en un planteo de prueba de hipótesis utilizando el modelo de la distribución normal. Veamos el siguiente ejemplo: Una investigadora analizó los promedios finales de la carrera de ingeniería de la UNC. Los docentes de esa Facultad emplean una evaluación en cada materia de puntuaciones de 1 al 10. Las estadísticas universitarias hasta el año 2012 indican que el promedio final de la carrera sigue una distribución normal con $M=6,33$ y $de=1,8$. Esta investigadora sospecha que durante los últimos años el promedio puede haber variado. Por ello, se decide hacer un estudio, registrando el promedio al final de la carrera de una muestra aleatoria de 200 alumnos.

La sospecha de la investigadora tiene que redactarse como una hipótesis de investigación y, además, plantearse en términos estadísticos.

H_1 : el promedio al final de la carrera de ingeniería de la UNC es diferente de los valores registrados hasta el año 2012.

H_0 : el promedio al final de la carrera de ingeniería de la UNC no es diferente de los valores registrados hasta el año 2012.

H_1 : $m_{2012} \neq m_{\text{actual}}$

H_0 : $m_{2012} = m_{\text{actual}}$

Nótese que en este ejemplo hemos introducido una nueva hipótesis que recibe el nombre de Hipótesis Nula, que registramos como H_0 y que es la negación de la Hipótesis de Investigación, que registramos como H_1 .

La pregunta es ¿por qué dos hipótesis? La respuesta es simple, porque la hipótesis que someteremos a prueba con los datos disponibles es la Hipótesis Nula. Hacemos esto puesto que H_0 es exacta y, por tanto, se ajusta a una distribución de probabilidad conocida. Para este ejemplo la distribución de probabilidad está especificada en la población de origen y es el modelo normal. Estadísticamente ambas hipótesis quedan redactadas de la siguiente manera.

$$H_0: \mu = 6,33$$

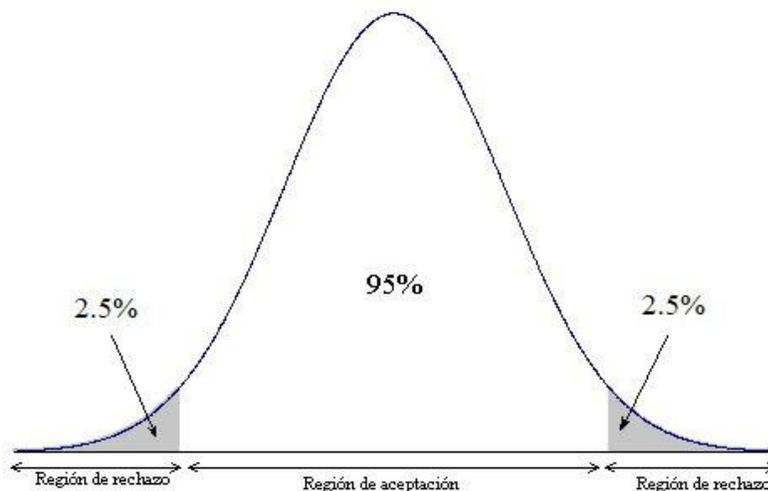
$$H_1: \mu \neq 6,33$$

Las hipótesis planteadas hasta aquí deben entenderse de la siguiente manera. La investigadora supone que el promedio de los estudiantes ya no es el mismo que aquel registrado hasta el año 2012, por lo tanto, los datos tomados de aquella población no representan adecuadamente a los estudiantes actuales. Sin embargo, no se puede decidir si el nuevo parámetro poblacional es mayor o menor que $\mu = 6,33$. Es por ello que la hipótesis de investigación solo especifica una diferencia. A este tipo de hipótesis se las denomina bidireccionales en tanto no anticipan la dirección de la diferencia. De esta manera, si el valor de μ encontrado en una nueva muestra de estudiantes tomada en el presente año está próxima al valor promedio 6,33, no hay razón para suponer que el valor del parámetro ha cambiado y por tanto se descarta la hipótesis de investigación. Al hacerlo se sostiene la hipótesis nula y en tal caso, la sospecha de la investigadora se resolvería concluyendo que no han existido cambios sustantivos en el promedio final de los egresados de la carrera de ingeniería y por tanto es lícito suponer que el parámetro para esa población es $\mu = 6,33$.

Hemos mencionado que no descartaremos la hipótesis nula siempre que el valor promedio al final de la carrera tomado de estudiantes en la actualidad, esté próximo a $\mu = 6,33$. Esto quiere decir que esperamos encontrar alguna fluctuación aleatoria en el valor de μ , pero que se encuentre dentro de ciertos límites de probabilidad. Estadísticamente, debemos establecer un límite dentro del cual aceptamos que el valor de μ sea diferente de 6,33 y que tal diferencia se deba a las fluctuaciones aleatorias del conjunto de datos. Una vez establecido este límite, el mismo servirá para determinar que las variaciones observadas no pueden ser explicadas por simples fluctuaciones aleatorias y por tanto no es factible sostener que el valor del parámetro es el que se especifica en la H_0 .

Como vimos anteriormente el modelo de la distribución normal nos ofrece un modelo para evaluar la validez de la H_0 . Según hemos visto por el teorema del límite

central, cualquier promedio de una muestra aleatoria extraída de una población tiende a ubicarse en el centro de esa distribución, de tal modo que: $\mu \sim \bar{x}$. Esto quiere decir que un conjunto de muestras tomadas aleatoriamente de esa población se ajusta a una distribución normal. Para determinar qué tan diferente es el valor de la muestra obtenida respecto del parámetro, solo basta comparar la desviación de μ respecto de \bar{x} . Para hacerlo utilizamos la puntuación Z correspondiente al modelo de la distribución normal estándar. Según hemos visto, el valor $Z=\pm 1,96$ se corresponde con un área que cubre el 95% de la distribución. De este modo podemos utilizar ese valor de referencia para determinar cuándo un promedio muestral se ha apartado lo suficiente del valor paramétrico para considerar que tal diferencia ya no puede atribuirse a simples variaciones aleatoria, sino que el valor del promedio muestral pertenece a una distribución con un valor paramétrico diferente. Bajo tales condiciones diremos que rechazaremos la H_0 planteada toda vez que encontremos un valor que sea mayor o menor que $Z=\pm 1,96$. Siendo así, hemos definido dos áreas dentro de la distribución normal, una en la que especifica los valores dentro de los cuales la H_0 debe ser aceptada, y otra en la que H_0 debe ser rechazada. Gráficamente esto se representa de la siguiente manera:



Tomando el valor crítico $Z=\pm 1,96$ tenemos que dentro de esos límites existe un 95% de posibilidades de que una muestra aleatoria de la población tenga un valor próximo al parámetro cuando es cierta H_0 . Rechazaremos la H_0 cuando su posibilidad de certeza sea baja, que según se muestra en la gráfica, ello ocurre cuando el valor muestral se aparta del centro de la distribución hacia alguno de los extremos. En tal

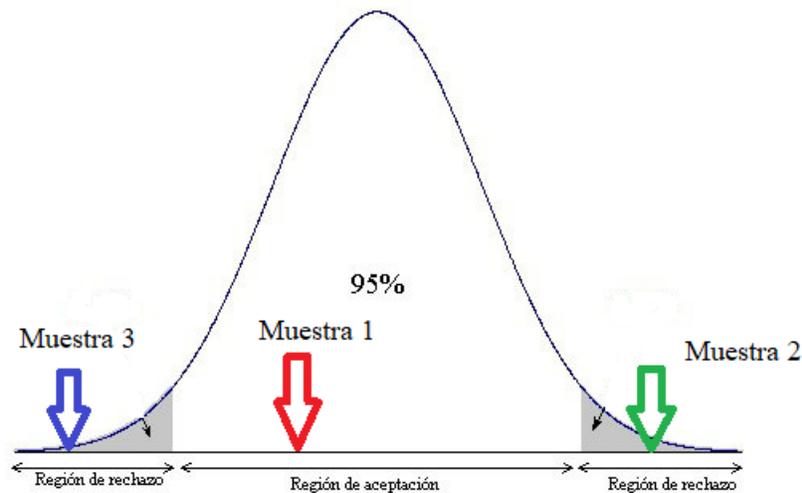
caso, se espera que H_0 sea cierta solo el 5% de las veces. En el caso de las hipótesis bidireccionales, al no especificar el sentido de la diferencia respecto al parámetro, se toman los dos extremos de la distribución, es por eso que la zona de rechazo de H_0 se encuentra repartida en un 2.5% por encima y por debajo de la distribución.

Vamos a introducir ahora un nuevo término para explicar cómo interpretar las zonas de aceptación y rechazo de H_0 . Estos términos son el Error tipo I y el Error tipo II. Cometemos un Error Tipo I toda vez que aceptamos una H_0 que es falsa, y cometemos el Error tipo II toda vez que rechazamos una H_0 que es verdadera. Ambos tipos de errores son complementarios y al ajustar uno de ellos también ajustamos el otro. Para simplificar la explicación, trabajaremos solo sobre el Error tipo I.

En el contexto de la prueba de hipótesis, queremos rechazar todas las H_0 que son falsas, pero eso es imposible debido a las fluctuaciones aleatorias que vimos anteriormente. Es por eso que el modelo probabilístico sobre el que trabajemos, nos proporciona valores críticos sobre los cuales establecer zonas de aceptación y rechazo de la H_0 . En la figura que se muestra a continuación, se han situado los valores promedios de tres muestras tomadas de una población con parámetro μ . Recordemos que la H_0 establece que $\mu \sim \bar{x}$

Por lo tanto, el promedio de la muestra 1 no es suficientemente diferente del parámetro como para rechazar la H_0 . En cambio, los valores de los promedios de las muestras 2 y 3 sí son lo suficientemente diferentes como para rechazar la H_0 .

Definimos un valor crítico para la aceptación o rechazo de la H_0 representándolo con la letra α (alfa). Si queremos tener una certeza del 95% de que al rechazar la H_0 estamos haciendo lo correcto, el valor de alfa es igual a 0.05. El valor de α se expresa como una probabilidad, pero es fácil comprender que esa probabilidad puede expresarse en porcentaje al multiplicarla por 100. Por tanto $\alpha=0.05$ es la probabilidad de cometer el Error Tipo I el 5% de las veces al rechazar una H_0 . Otra manera de decirlos es, que $\alpha=0.05$ nos da una certeza de que el 95% de las veces no rechazaremos H_0 cuando es cierta. Por lo tanto, minimizamos el Error Tipo I, si rechazamos H_0 cuando el valor promedio obtenido fuera el de las muestras 2 y 3. En el caso de la muestra 1 lo mejor es no rechazar la H_0 .



Volvamos ahora al ejemplo que planteamos al comienzo. Una investigadora analizó los promedios finales de la carrera de ingeniería de la UNC. Los docentes de esa Facultad emplean una evaluación en cada materia de puntuaciones de 1 al 10. Las estadísticas universitarias hasta el año 2012 indican que el promedio final de la carrera sigue una distribución normal con $M=6,33$ y $de=1,8$. Esta investigadora sospecha que durante los últimos años el promedio puede haber variado. Por ello, se decide hacer un estudio, registrando el promedio al final de la carrera de una muestra aleatoria de 200 alumnos. La hipótesis planteada y su modelo estadístico es el siguiente:

H_1 : el promedio al final de la carrera de ingeniería de la UNC es diferente de los valores registrados hasta el año 2012.

H_0 : el promedio al final de la carrera de ingeniería de la UNC no es diferente de los valores registrados hasta el año 2012.

$H_1: \bar{x}_{2012} \neq \bar{x}_{\text{actual}} \quad H_1: \mu \neq 6.33$

$H_0: \bar{x}_{2012} = \bar{x}_{\text{actual}} \quad H_0: \mu = 6.33$

Dado que la investigadora no predice una diferencia en ninguna dirección H_1 es Bidireccional, por lo tanto, estaremos usando los dos extremos del modelo de distribución normal estándar. Es decir, la prueba será de dos extremos.

Ahora debemos establecer el valor crítico de rechazo de H_0 , es decir, debemos darle un valor a α . Para continuar con lo argumentado hasta aquí, establecemos $\alpha=0.05$. Siendo este el valor establecido, sabemos ahora que los valores de Z críticos que establecen las zonas de aceptación y rechazo de H_0 corresponden a los valores $Z \pm 1.96$.

Para finalizar el ejemplo, supongamos que los resultados del estudio realizado por la investigadora, determinaron que el promedio registrado en la muestra es de $\bar{x}=5.75$ con una $de= 1,6$.

Para obtener el valor Z observado que compararemos con el valor Z crítico, necesitamos obtener el error estándar de la media. La distribución del error estándar de la media es el mismo que utilizamos para la estimación de parámetros, salvo que el denominador utiliza el valor del tamaño muestral.

Continuando con el ejemplo, diremos que para el estudio que realiza la investigadora, se utiliza una muestra aleatoria de tamaño $n=200$. De los cual resulta que el error estándar de la media para ese tamaño muestral es de 0,113. La siguiente ecuación resume lo dicho para la distribución de datos del ejemplo:

$$EE_{\bar{x}} = \frac{DE}{\sqrt{n}} = \frac{1.6}{\sqrt{199}} = 0,113$$

Sabemos que el valor Z crítico para este ejemplo se estableció en $Z\pm 1.96$. Debemos, encontrar el valor de Z observado para realizar la comparación. Sabemos que para encontrar una puntuación Z debemos proceder con la siguiente ecuación:

$$Z = \frac{x - \bar{x}}{DE}$$

Para la prueba de hipótesis, debemos reemplazar en la ecuación el promedio muestral por el parámetro, y la desviación estándar por el $EE_{\bar{x}}$, de modo que el valor Z observado se resuelve de la siguiente manera.

$$Z_{observado} = \frac{\bar{x} - \mu}{EE_{\bar{x}}}$$

Asumiendo que, en la investigación realizada, se encontró que el promedio final de los estudiantes de ingeniería es de 6,86, el cálculo del valor Z observado da como resultado:

$$Z_{\text{observado}} = \frac{6.86 - 6.33}{0.113} = 4.6$$

La zona de rechazo de H_0 se encuentra en algún valor entre $Z \pm 1.96$. Encontramos que el valor de Z observado supera largamente el valor de Z crítico establecido para rechazar H_0 . Asumimos entonces que, siendo $Z_{\text{observado}} > Z_{\text{crítico}}$ para el valor de $\alpha = 0.05$, se rechaza H_0 . Por lo tanto, la hipótesis de la investigadora es correcta al suponer que el parámetro ha cambiado.

Habiendo rechazado la H_0 nos queda como válida H_1 , y en este caso corresponde recalibrar el parámetro μ para el promedio de los estudiantes de ingeniería al final de la carrera realizando para ello una nueva estimación del mismo.

La hipótesis sobre la que realizamos todo el planteo es de tipo bidireccional en tanto no anticipamos que el nuevo parámetro estuviera por encima o por debajo del parámetro actual. Sin embargo, se podría haber previsto esa diferencia tal como se hizo con la hipótesis planteada al principio de este apartado sobre los operativos de evaluación.

Retomemos el ejemplo sobre el que estamos trabajando y reformulemos H_1 . Una investigadora analizó los promedios finales de la carrera de ingeniería de la UNC. Los docentes de esa Facultad emplean una evaluación en cada materia de puntuaciones de 1 al 10. Las estadísticas universitarias hasta el año 2012 indican que el promedio final de la carrera sigue una distribución normal con $M = 6,33$ y $de = 1,8$. Esta investigadora sospecha que durante los últimos años el promedio ha variado y en la actualidad es más alto. Por ello, se decide hacer un estudio, registrando el promedio al final de la carrera de una muestra aleatoria de 200 alumnos.

El modelo estadístico de la hipótesis planteada es el siguiente:

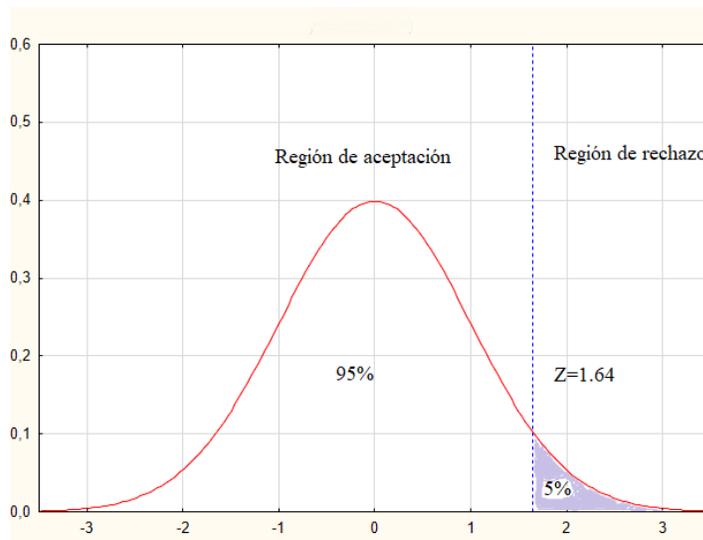
$$H_1: \bar{x}_{2012} < \bar{x}_{\text{actual}} \quad H_1: \mu > 6,33$$

$$H_0: \bar{x}_{2012} = \bar{x}_{\text{actual}} \quad H_0: \mu = 6,33$$

Vemos que el planteo supone ahora que la diferencia entre el parámetro establecido a partir de los datos de 2012 es menor al actual. Por su parte H_0 sigue siendo la misma en tanto se trata de negar a la H_1 .

De todo este planteo, lo único que se modifica es la zona en donde estableceremos el punto de aceptación o rechazo de la hipótesis nula. Si examinamos la H_1 vemos que es unidireccional en tanto determina que el nuevo parámetro es mayor, de modo que el valor del punto de corte Z debe situarse sobre el extremo superior de la curva. Este nuevo punto de corte debe establecer las zonas de

aceptación y rechazo de H_0 . Si no cambiamos el ejemplo más que en el punto mencionado, el nuevo valor de Z crítico se establece es $Z_{\text{crítico}}=1,64$. Gráficamente ello se representa como se muestra a continuación.



Vemos que toda la zona de aceptación queda delimitada por el área de la izquierda de la curva y el valor Z crítico se sitúa en el extremo derecho de la distribución, a partir del valor 1,64.

En el ejemplo, habíamos encontrado que el valor de Z observado= 4.6 que es mayor que 1.64. De este modo, bajo este nuevo planteo de hipótesis también rechazaríamos H_0 .

Finalmente, cabe destacar que, para las hipótesis unidireccionales, se utilizan los extremos de la distribución normal, de acuerdo al modo en que se plantea la hipótesis de investigación, es decir, de acuerdo a la manera en que anticipamos en que ocurrirá la diferencia. Solo para las hipótesis bidireccionales utilizamos ambos extremos de la distribución simultáneamente.

Capítulo 8

Prueba de hipótesis y el modelo estadístico χ^2

Supongamos que un investigador dispone de una muestra de 100 profesores, algunos de ellos dictan clases en el último ciclo de primaria y otros en el primer año de Secundaria. A todos ellos se les hizo la siguiente pregunta: ¿Utiliza usted la evaluación como una instancia de aprendizaje? La respuesta requerida se clasificó como una variable dicotómica nominal: SI - NO. El investigador sospecha que aquellos profesores que dictan clases en los cursos de primaria tienden a usar las evaluaciones como instancias de aprendizaje, lo cual no ocurre con los profesores que dictan clases en la Secundaria. Así, el investigador plantea la siguiente hipótesis:

H₁: Existe relación entre el nivel de escolaridad donde el profesor dicta sus clases y la utilización de la evaluación como instancia de aprendizaje.

Nótese entonces que la hipótesis de investigación es una formalización de lo que él intuye que ocurre con las evaluaciones en los dos ciclos de escolarización. Como ya sabemos, la hipótesis de investigación en términos conceptuales siempre se acompaña de la hipótesis nula (H₀). Esta última toma la forma de una negación de la hipótesis de investigación; en otras palabras, es su contrario. Por lo tanto, la misma se formaliza como:

H₀: No existe relación entre el nivel de escolaridad donde el profesor dicta sus clases y la utilización de la evaluación como instancia de aprendizaje.

Si hay asociación entre las variables deberíamos esperar que las proporciones de aquellos profesores que responden afirmativamente (o negativamente) se concentren preferentemente en una categoría de nivel de escolaridad. Al contrario, si no existe asociación entre las variables deberíamos esperar que la respuesta de los profesores estuviera repartida proporcionalmente en cada una de las categorías de nivel de escolaridad. En otras palabras, si no existe asociación entre las variables, las frecuencias en la categorías SI y NO, estarían proporcionalmente repartidas entre los docentes que enseñan en primaria y en Secundaria. Al cruzar las variables nivel de escolaridad y respuestas de los profesores, es posible construir una tabla de contingencia de cuatro celdas, tal como se muestra a continuación.

	¿Utiliza la evaluación como instancia de aprendizaje?	
Nivel donde Dicta clases	SI	NO
EGB		
Secundaria		

Las celdas están allí donde se entrecruzan los niveles de las variables. Bajo el supuesto de la hipótesis nula, las frecuencias de las celdas deberían repartirse proporcionalmente de acuerdo a las frecuencias marginales. Cabe destacar que la sospecha del investigador es que los profesores del último ciclo de primaria son los que utilizan la evaluación como instancia de aprendizaje, pero esta situación no está expresada en la hipótesis planteada. Es decir, si la hipótesis debiera responder estrictamente a la sospecha del investigador, esta debería haberse planteado de modo unidireccional. Pero, por tratarse solo de una sospecha, nuestro investigador solo plantea una hipótesis bidireccional.

Siendo P_{si} la proporción de profesores que responden afirmativamente y P_{no} la proporción de profesores que responden negativamente a la pregunta del investigador, las hipótesis toman la siguiente forma estadística de acuerdo a la categoría: ciclo en el que dicta clase el profesor:

$$H_1: p_{si} \neq p_{no}: \chi^2 > 0$$

$$H_0: p_{si} = p_{no}: \chi^2 = 0$$

Cálculo del estadístico χ^2 a partir de la tabla de contingencia

Siguiendo con el ejemplo planteado, supondremos que el investigador cuenta con las respuestas dadas por los profesores, y vuelca los datos en una tabla de contingencia, colocando en cada una de las entradas las variables de interés, y arreglando las frecuencias observadas en cada una de las casillas. La frecuencia total que resulta de sumar las frecuencias en cada una de las casillas que forman las filas y las columnas se denominan marginales de fila y de columna respectivamente. De este modo, los datos se presentarían tal como lo muestra la siguiente tabla.

	¿Utiliza la evaluación como instancia de aprendizaje?		
Nivel donde dicta clases	SI	NO	
Primaria	33	19	52
Secundaria	27	21	48
	60	40	100

Esta tabla muestra que a la pregunta respondieron 52 profesores de primaria y 48 de secundaria; además se observa que 60 de estas preguntas fueron afirmativas y 40 negativas. Estos son los marginales por fila y columna de la tabla. En cada una de las celdas tenemos las proporciones de respuestas afirmativas y negativas de los profesores de primaria y Secundaria. Estas son las frecuencias observadas, es decir, aquellas que se obtienen empíricamente del conteo de las unidades de análisis que caen en el cruce de cada una de las categorías de las variables.

Anteriormente se dijo que bajo la hipótesis nula debería esperarse una distribución proporcional de las frecuencias en cada una de las celdas. Conociendo las frecuencias marginales es posible calcular las frecuencias que cabría esperar en cada una de las celdas, si efectivamente no existiera asociación entre las variables. Estas frecuencias se llaman frecuencias esperadas porque son aquellas frecuencias teóricas que cabría esperar, de ser cierta la hipótesis nula. Su cálculo es muy sencillo, se trata de multiplicar el marginal de fila y columna correspondiente a cada celda y dividirlo por el total de casos. Veamos en detalle cómo se obtienen:

Frecuencias esperadas para:

- a) Profesores de primaria que respondieron afirmativamente: $60 \cdot 52 / 100 = 31.2$
- b) Profesores de primaria que respondieron negativamente: $40 \cdot 52 / 100 = 20.8$
- c) Profesores de Secundaria que respondieron afirmativamente: $60 \cdot 48 / 100 = 28.8$
- d) Profesores de Secundaria que respondieron negativamente: $40 \cdot 48 / 100 = 19.2$

Entonces, las frecuencias observadas son las que efectivamente se cuentan como cantidad de casos en cada una de las celdas que surgen del cruce entre las categorías de las variables. Las frecuencias esperadas son teóricas y se obtienen de la operación de multiplicar las frecuencias marginales correspondientes a cada celda y dividirla por el total de casos. Ahora, se puede arreglar las frecuencias observadas y esperadas en la misma tabla y compararlas. Cabe recordar en este punto que las frecuencias esperadas son las que se hubieran obtenido de ser cierta la hipótesis nula, y puesto que ésta representa la situación en donde no hay asociación entre las variables, la proporción de las frecuencias esperadas es la misma que la proporción en los marginales. Así por

ejemplo, la proporción de profesores de primaria sobre profesores de Secundaria es de $52/48= 1,083$, la misma proporción entre los profesores que respondieron afirmativamente es de $31,2/28,8= 1,083$, y la proporción entre los profesores que respondieron negativamente es de $20,8/19,2= 1,083$ (el mismo cálculo puede efectuarse con los marginales de columnas). Entonces las frecuencias esperadas representan una distribución proporcional de casos en cada una de las celdas, dados los marginales de tabla.

Comparación de frecuencias observadas y esperadas

Nivel donde dicta clases	¿Utiliza la evaluación como instancia de aprendizaje?		
	SI	NO	
EGB	33 (31.2)	19 (20.8)	52
Secundaria	27 (28.8)	21 (19.2)	48
	60	40	100

Las frecuencias esperadas se muestran entre paréntesis

Según lo que se ha explicado, si las dos variables estuvieran relacionadas entre sí, la distribución de frecuencias no debería ser proporcional, dado que una variable ejercería alguna “acción” sobre la otra. En este ejemplo, se supone que los profesores que enseñan en distintos niveles, toman a la evaluación de manera diferente. Por lo tanto, si las frecuencias observadas están próximas o son iguales a las frecuencias esperadas, estaríamos en la situación planteada por la hipótesis nula. Al contrario, si las frecuencias observadas se apartan notablemente de las esperadas será necesario rechazar la hipótesis nula. El estadístico χ^2 toma en consideración la distancia entre las frecuencias observadas y las esperadas, y sirve para estimar la probabilidad de que tales diferencias sean debidas al azar. Veamos primero la forma de cálculo de este estadístico y luego retomaremos este punto.

La fórmula de cálculo de χ^2 puede representarse de la siguiente manera:

$$\chi^2 = \frac{\sum (f_o - f_e)^2}{f_e}$$

El estadístico χ^2 es la sumatoria de las desviaciones cuadráticas entre las frecuencias observadas y esperadas, sobre las frecuencias esperadas. Dado el caso en que las frecuencias observadas sean las mismas que las frecuencias esperadas, el numerador de la fracción es igual a cero, por lo tanto el estadístico también es igual a cero. Esta es la situación exacta que plantea la hipótesis nula. En otro caso tendremos que las frecuencias observadas serán distintas de las frecuencias esperadas y el estadístico se

apartará progresivamente de cero. Por ende, cuanto mayor sea la diferencia entre las frecuencias observadas y las esperadas, más alejados estaremos de la situación planteada por la hipótesis nula. En el ejemplo el estadístico χ^2 se calcula de la siguiente manera:

$$\chi^2 = (33-31.2)^2/31.2 + (19-20.8)^2/20.8 + (27-28.8)^2/28.8 + (21-19.2)^2/19.2 = \mathbf{0.538}$$

El valor del estadístico χ^2 debería servirnos para decidir si rechazamos o no la hipótesis nula. Para ello debemos tomar en cuenta la distribución del estadístico, y el valor de probabilidad asignado al Error Tipo I.

Reglas de decisión basadas en χ^2 y errores de decisión.

Cuando se plantea una hipótesis de investigación, tenemos también una hipótesis nula. En el ejemplo que estamos desarrollando, la hipótesis nula es analizada bajo el modelo del estadístico χ^2 para decidir si es posible rechazarla. Nótese que anteriormente planteamos que la hipótesis nula es una negación de la hipótesis de investigación, de modo que si es posible rechazar la hipótesis nula, la hipótesis de investigación es la alternativa válida. En caso que no sea posible rechazar la hipótesis nula, debemos descartar la hipótesis de investigación. En otras palabras, al rechazar la hipótesis nula se asume que efectivamente existe asociación entre las variables y que tal asociación no es producto del azar.

Sabemos que al tomar la decisión de rechazar la hipótesis nula se puede cometer un error, rechazarla cuando esta es verdadera. Si procedemos de tal manera habremos cometido un error tipo I. En cambio, si aceptamos una hipótesis nula cuando es falsa, estamos cometiendo un error tipo II. No es posible disminuir al mínimo ambos tipos de errores al mismo tiempo, pero en un planteo de hipótesis es factible establecer un equilibrio entre ambos. Vimos que una manera de hacerlo es establecer una condición donde la probabilidad de cometer el error tipo I sea la menor posible. Convencionalmente, en una investigación la probabilidad de cometer el error tipo I se establece a priori y se define como α (alfa), este valor se determina por convención en $\alpha=0.05$. Mediante el valor de α podemos establecer un valor crítico para χ^2 , con el cual compararemos el valor χ^2 efectivamente observado. La regla de decisión quedará expresada de la siguiente forma:

Sea $\alpha=0.05$, se rechaza la hipótesis nula siempre y cuando:

$$\chi^2 \text{ observado} \geq \chi^2 \text{ crítico}$$

Es decir, si el valor de χ^2 observado es igual o mayor al valor de χ^2 crítico, rechazamos la hipótesis nula que establece que las variables no están asociadas. Si, al contrario, el valor de χ^2 observado no supera el valor de χ^2 crítico, no hay evidencia suficiente para rechazar la hipótesis nula.

Ya se dijo que el valor de χ^2 crítico se obtiene al establecer el valor de α . Ahora bien, hay que tener en cuenta que la distribución de este estadístico conforma una familia de curvas, que depende del tamaño de la tabla de contingencia. Es decir, el valor de χ^2 crítico varía según los grados de libertad de la tabla. Aunque esto parece engorroso, en realidad es un concepto muy sencillo, digamos que para conocer cuántos grados de libertad tiene una tabla debemos resolver la siguiente fórmula:

$$gl = (f-1)*(c-1)$$

Entonces, los grados de libertad de la tabla se definen como la cantidad de filas menos 1, multiplicado por la cantidad de columnas menos 1. En la tabla de nuestro ejemplo tenemos 2 filas, una que corresponde a profesores de primaria y la otra corresponde a profesores de secundaria. También tenemos dos columnas, una que corresponde a todos los que respondieron SI y la otra la que corresponde a todos los que respondieron NO. Al aplicar la fórmula tendremos que los grados de libertad de la tabla resultan, $gl = (2-1)*(2-1) = 1$. El valor de χ^2 crítico que estamos buscando para resolver nuestro ejemplo, es aquel que corresponde a una tabla de contingencia con 1 grado de libertad, para un valor de $\alpha = 0.05$.

Existen tablas donde se ha tabulado el valor de χ^2 crítico para diferentes grados de libertad, y para distintos valores de α . En nuestro caso el valor buscado de χ^2 crítico es de 3,841. Al aplicar la regla de decisión planteada tenemos que χ^2 observado $<$ χ^2 crítico esto es, $0,538 < 3,841$ y por lo tanto no es posible rechazar la hipótesis nula.

De acuerdo a lo observado en la tabla 9, tenemos que no existe asociación entre el nivel en que enseña el profesor y el hecho de utilizar la evaluación como una instancia de aprendizaje. Ahora el investigador sabe que no existe asociación entre esas variables, y que su sospecha posiblemente se debió a un sesgo en sus observaciones.

En el siguiente cuadro se ofrecen valores críticos de χ^2 que corresponde a valores de $\alpha=0,05$ y $\alpha=0,01$ y diferentes dimensiones de la tabla de contingencia.

Gl	ET I=1%	ET I= 5%
1	6,63	3,48
2	9,21	5,99
3	11,34	7,81
4	13,27	9,48

El grado de asociación entre las variables

En el ejemplo que hemos dado, nos encontramos con que no ha sido rechazada la hipótesis nula y por tanto se asume que no hay asociación entre las variables. Pero en el caso en que se haya rechazado la hipótesis nula, asumiremos que las variables están relacionadas y será necesario determinar el grado de asociación entre ellas. Basados en la distribución χ^2 existen varias alternativas para calcular la asociación entre las variables. En los apartados siguientes veremos dos estadísticos que determinan esa asociación y que son los más comúnmente usados.

Coefficiente de contingencia C de Pearson

Cuando se encuentra correlación entre las variables e interesa conocer el grado de asociación entre ellas, se emplea el coeficiente de contingencia C de Pearson, el cual se define como: la raíz cuadrada del cociente entre el valor de χ^2 y el valor de χ^2 más el número de casos N.

$$C = \sqrt{\frac{\chi^2}{\chi^2 + N}}$$

El coeficiente de contingencia C de Pearson, (que no debe ser confundido con el coeficiente de correlación r de Pearson), estima la asociación entre variables en una tabla de contingencia. Cuando estas son independientes, el coeficiente es igual a cero, y se hace mayor a medida que aumenta la dependencia de las variables. El valor máximo del coeficiente depende del número de filas y columnas en la tabla de contingencia, pero es siempre menor que la unidad.

El hecho que este coeficiente no tenga un valor máximo establecido para cualquier tabla, dificulta su interpretación, pero al ser constante permite la comparación entre resultados obtenidos con tablas de igual tamaño.

Coeficiente V de Cramer

Otro coeficiente de asociación que puede utilizarse para estimar el grado de asociación entre los atributos es el coeficiente V de Cramer, cuya fórmula es la siguiente:

$$V = \sqrt{\frac{\chi^2}{\min(h - 1; c - 1) * N}}$$

En la fórmula la expresión $\min. (h-1; c-1)$ alude al número de hileras o columnas menos uno. Es decir, en el numerador entrará el menor valor de los dos. Entonces V se define como: la raíz cuadrada del cociente entre el valor χ^2 observado y el número de columnas o filas menos uno (dependiendo de cuál sea el menor), multiplicado por el número de casos. Este coeficiente puede aplicarse a tablas de cualquier tamaño y alcanza un máximo de uno.

Ejemplo de aplicación

Vamos a presentar un ejercicio de aplicación de la prueba χ^2 para el caso de una tabla de contingencia mayor a 2x2. Un investigador realiza una encuesta a un grupo de docentes de primaria, que enseñan en escuelas privadas y públicas, sobre la importancia de trabajar en equipos interdisciplinarios cuando en el aula existen niños con problemas de aprendizaje. El investigador desea comprobar si existe asociación entre los atributos, tipo de escuela y valoración del trabajo interdisciplinario. Plantea la hipótesis de que los docentes darían distinta valoración al trabajo en equipo, de acuerdo a la gestión de la escuela en que trabajan. Para comprobar la existencia de asociación entre las variables utiliza el estadístico χ^2 . De acuerdo a las hipótesis esbozadas se tiene que el planteo estadístico correspondiente es:

$$H_1: \chi^2 > 0$$

$$H_0: \chi^2 = 0$$

La tabla recoge la distribución de frecuencias observadas que resulta de la combinación de las variables gestión de la escuela y valoración del trabajo interdisciplinario.

Frecuencias Observadas

Gestión de la escuela	Valoración del trabajo interdisciplinario		
	Es muy importante	Es importante pero no necesario	No es tan importante
Privadas	102	106	21
Públicas	195	92	89

Corresponde ahora construir la tabla que hubiéramos encontrado de ser verdadera la hipótesis nula. Dicha tabla se construye con las frecuencias esperadas y refleja la situación en que el estadístico χ^2 es igual a 0.

Frecuencias Esperadas

Gestión de la escuela	Valoración del trabajo interdisciplinario		
	Es muy importante	Es importante pero no necesario	No es tan importante
Privadas	112,42	74,94	41,63
Públicas	184,58	123,05	68,36

La diferencia entre la distribución de frecuencias observadas y esperadas, especificará la magnitud del estadístico χ^2 . Para ello debemos calcular la sumatoria de esas diferencias, tal como se muestra a continuación:

1	$(102-112.42)^2/112.42=0.965$
2	$(195-184.58)^2/184.58=0.588$
3	$(106-74.94)^2/74.94=12.87$
4	$(92-123.05)^2/123.05=7.83$
5	$(21-41.63)^2/41.63=10.22$
6	$(89-68.36)^2/68.36=6.23$
Σ	38.703

Hemos obtenido el valor de χ^2 observado y corresponde compararlo con el valor χ^2 crítico, para tomar la decisión de rechazar o no la hipótesis nula. El valor de χ^2 crítico para un error tipo I igual a $\alpha=0.05$ y una tabla con 2 grados de libertad es de 5,99. Entonces tenemos que el valor de χ^2 observado es mayor que el valor de χ^2 crítico; de acuerdo a la regla de decisión corresponde rechazar la hipótesis nula.

Aceptando que los atributos de las variables están relacionados, corresponde determinar el grado de asociación entre las variables. Para ello aplicamos el coeficiente V de Cramer.

$$V = \sqrt{\frac{38,703}{605}} = 0,252$$

Encontramos que las variables están asociadas, pero la magnitud de dicha asociación es moderada a débil. ¿Cuál es la manera de interpretar esa asociación? Si repasamos la tabla vemos que la mayor diferencia se produce en el nivel de la variable valoración del trabajo interdisciplinario para escuelas de gestión privadas, en menor medida se registran discrepancias en las escuelas de gestión públicas. Por lo tanto, es factible concluir que la asociación entre las variables está dada por la inclinación de los docentes de escuelas privadas a opinar que el trabajo interdisciplinario es importante pero no necesario y los docentes de escuelas públicas a opinar que no es tan importante. Aunque la asociación resulta de moderada a débil, tal inclinación resulta atendible en virtud de que el hemos rechazado la hipótesis de no asociación entre las variables.

Las frecuencias esperadas como eventos independientes

Hasta aquí hemos señalado que las frecuencias esperadas son las que ocurrirían de ser cierta la hipótesis nula y además, es la distribución de frecuencias que hace que el estadístico χ^2 sea igual a cero. Esto, equivale a decir que cualquier evento que conste en las columnas será independiente de cualquier otro evento que conste en las filas. Volvamos por un momento a la tabla precedente. Si no hay asociación esto quiere decir que los profesores que dictan clase en escuelas privadas pueden valorar de cualquier manera el trabajo interdisciplinario. Un razonamiento similar puede hacerse si se parte desde las filas. De modo que la frecuencia esperada, resulta de tomar como independientes las filas y las columnas y obteniendo de ello una probabilidad. La regla de la multiplicación de probabilidades de eventos independientes sirve para entender por qué las frecuencias esperadas se calculan de esta manera.

Por ejemplo, si quisiéramos conocer la probabilidad de obtener un resultado particular cuando se lanzan al aire dos veces la misma moneda, se aplica la regla multiplicativa de las probabilidades independientes, dado que suponemos que un evento (el primer lanzamiento de la moneda), no influye en el segundo evento (el segundo lanzamiento de la moneda). Llamamos A al primer evento y B al segundo; ambos son independientes. Entonces, la probabilidad de un resultado AB se expresa como $p(AB) = p(A) \times p(B)$. Para el caso de la moneda, en la cual podemos obtener solo dos resultados posibles, la probabilidad de A es igual a $1/2 = 0.5$ (la probabilidad de B será la misma que A). Según lo expresado en la fórmula, la probabilidad de dos veces cara será $p^2 \text{ caras} = (0.5) \times (0.5) = 0.25$.

Si volvemos por un momento a nuestro ejemplo y considerando cierta la hipótesis nula: el tipo de gestión de la escuela donde dicta clases el profesor no se relaciona con la valoración del trabajo interdisciplinario, se puede buscar la

probabilidad de los eventos combinados, siempre considerando que estos eventos son independientes. Se puede calcular por un lado, la probabilidad de obtener de la muestra de profesores, uno que pertenece al tipo de gestión privado. Para ello consideramos a todos los profesores en ese nivel de la variable que suman 229 profesores (este es el marginal de fila) y lo dividimos por el total $229/605=0,378$. Asimismo, podemos calcular la probabilidad de obtener de todos los profesores, uno que considere el trabajo interdisciplinario como muy importante. Para ello debemos considerar a todos los profesores de la primer columna que suman 297 (este es el marginal de columna) y dividimos por el total $297/605=0,49$. La probabilidad de ambos eventos combinados resulta de la multiplicación de los eventos independientes, esto es: $0,378 \times 0,49 = 0,18522$. El valor obtenido es la frecuencia relativa o valor de probabilidad, pero puesto que en la tabla se tienen las frecuencias absolutas se debe transformar la frecuencia obtenida. Por lo que ya sabemos, para pasar de una frecuencia relativa a una absoluta, hay que multiplicarla por N, entonces $0,185 \times 605 = 112,0581$, es la frecuencia esperada buscada para la intersección de las variables Profesor de escuela privada que valora el trabajo interdisciplinario como muy importante (la diferencia en decimales entre el valor obtenido y el de tabla se debe al redondeo). El proceso de multiplicar marginal de fila por marginal de columna y dividirlo por el total de casos es una simplificación de lo anteriormente dicho.

Resumiendo, las frecuencias esperadas representan la probabilidad de un evento combinado, cuando ambos eventos son independientes, y la independencia de los eventos es el postulado de la hipótesis nula. Asimismo, el estadístico χ^2 es una comparación entre las frecuencias observadas y las esperadas, bajo la hipótesis de la independencia de los eventos.

Los grados de libertad en tablas de contingencia

Sabemos que la frecuencia marginal de una fila se construye a partir de la suma de las frecuencias observadas en cada una de las celdas. Así, calculamos la frecuencia marginal de la fila correspondiente a profesores de EGB, sumando las frecuencias de celdas que son 33 y 19, por tanto esta frecuencia marginal es igual a 52 (ver tabla). Sabiendo ahora que el total de esa suma es 52, ¿cuántos valores pueden variar libremente si se necesitan dos sumandos para ese resultado? Para exponerlo gráficamente, estaríamos en esta situación:

$$_ + _ = 52$$

Si se completa el primer espacio en blanco, el segundo espacio queda determinado por defecto, así si en el primer espacio se coloca el número 1, en el segundo espacio se

debe colocar el número 51 para que la suma sea igual a 52. Si en vez de colocar el número 1 se coloca el número 2, en el segundo casillero se debe colocar el número 50 para que la suma siga siendo igual a 52. Es decir, cuando tengo un resultado dado, que se obtiene por la suma de dos miembros, solo uno puede variar libremente. En una tabla de contingencia de 2x2 cada uno de los marginales es el resultado de una suma de dos números, por lo tanto solo un número en una de las cuatro celdas puede variar libremente. Sigamos el procedimiento esbozado más arriba. La menor frecuencia marginal de la tabla es igual a 40 (que es la suma de las celdas de los profesores de EGB y Secundaria que respondieron negativamente). Sabiendo que ningún miembro de la suma puede ser mayor que el total, diremos que el límite para repartir libremente los valores dentro de la tabla es de 40. Construyamos entonces una primera tabla asignando el valor 1 a la casilla correspondiente a profesores de primaria que responden negativamente:

	Utiliza la evaluación como instancia de aprendizaje?		
Dicta clases	SI	NO	
Primaria	51	1	52
Secundaria	9	39	48
	60	40	100

Si determinamos ese valor de celda arbitrariamente, el resto de los valores queda establecido, ya que no hay otros valores posibles que den como resultado los marginales de tabla. Sigamos con ese procedimiento y asignemos el valor 2 a esa casilla. Tal como muestra la tabla, al asignar el número 2 los únicos valores de las restantes son los siguientes.

	Utiliza la evaluación como instancia de aprendizaje?		
Nivel donde Dicta clases	SI	NO	
Primaria	50	2	52
Secundaria	10	38	48
	60	40	100

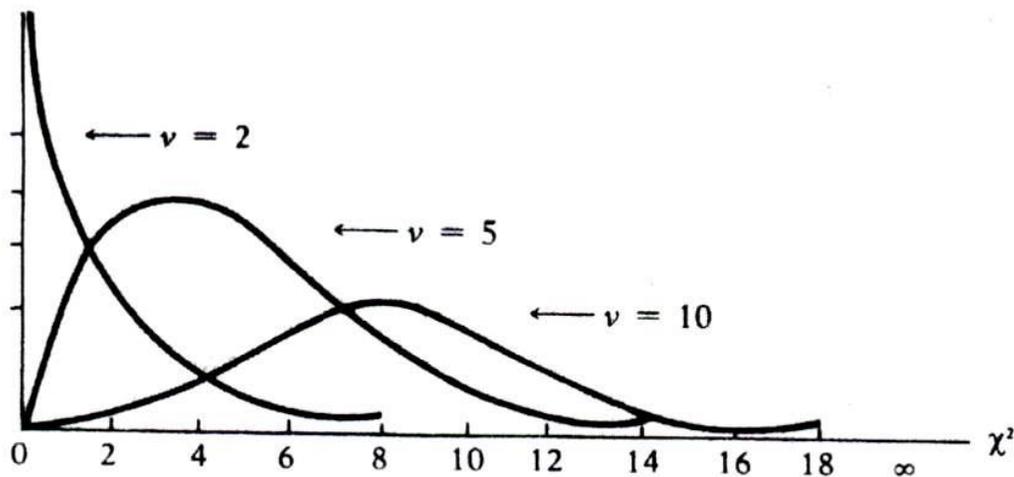
Como ya vimos, una tabla 2x2 esta tendrá 1 solo grado de libertad. Ahora bien, resulta que si la tabla tiene más filas y más columnas, tendrá más grados de libertad.

Dados los marginales de la tabla, ¿cuántas tablas de contingencia diferentes pueden construirse? De hecho sabemos que se pueden construir al menos 40 tablas distintas. Determinar la cantidad de tablas es un cálculo complicado, pero si sabemos que el menor valor de frecuencia marginal de la tabla es 40, y si además sabemos que ese valor se obtiene de la suma de dos celdas, y que determinado el valor de una de

esas celdas, los valores de las restantes celdas quedan determinados por defecto, se concluye que se pueden construir varias tablas distintas con los marginales dados. Ahora bien, de construir efectivamente esas tablas, sabríamos que algunas de ellas representarían una distribución proporcional de frecuencias. Es decir, algunas tablas estarían en consonancia con lo planteado por la hipótesis nula, mientras que otras presentarían valores en las celdas que harían insostenible esa hipótesis. Entonces, ¿qué probabilidad existe de haber encontrado una distribución de frecuencias como la observada en la tabla, bajo el supuesto que la hipótesis nula es verdadera? En tal caso, revisaríamos la probabilidad dada por la prueba χ^2 la cual establece que para una hipótesis bidireccional como la planteada, la probabilidad de ocurrencia de una tabla con una distribución de frecuencia observada, ocurriría con una probabilidad mayor al 50% (concretamente el 54,15%). Tal probabilidad de ocurrencia es alta para restarle credibilidad a la hipótesis nula.

Distribución χ^2 .

Como ya se dijo antes, la distribución χ^2 es una distribución conocida conformada por una familia de curvas; estas diferentes curvas se basan en los grados de libertad, que como vimos, no es otra cosa que la cantidad de celdas que están libres de variar, cuando se han establecido los marginales de la tabla. Como se aprecia en la siguiente figura, la distribución es asimétrica y se aproxima a una distribución normal cuando los grados de libertad son mayores que 10. En todos los casos la distribución muestra la probabilidad de ocurrencia de un valor χ^2 , bajo el supuesto de que la hipótesis nula es verdadera. Se observa entonces que si el valor está próximo a 0, la probabilidad de que la hipótesis nula sea verdadera es alta, y al contrario, a medida que aumenta el valor de χ^2 esa probabilidad se hace cada vez menor. El valor de χ^2 crítico se sitúa sobre el eje de las abscisas y resulta el punto de comparación con el valor de χ^2 observado. En este caso es posible utilizar la regla de decisión esbozada anteriormente o bien, determinar el valor exacto de ocurrencia del estadístico χ^2 obtenido.



Recomendaciones para el uso de la prueba χ^2 .

La prueba requiere que las frecuencias esperadas en cada celda no sean demasiado pequeñas, ya que si este es el caso, la prueba pierde potencia. Las recomendaciones más citadas para la correcta aplicación son:

Tablas de 2×2 ; $gl=1$: cuando n es igual o menor que 20, se recomienda el uso de la prueba exacta de Fisher. Cuando n es mayor que 20 y menor que 40, la prueba χ^2 puede utilizarse si las frecuencias esperadas alcanzan un valor no menor que 5. Si existen frecuencias esperadas menores que 5, conviene utilizar la prueba exacta de Fisher.

Tablas de $n \times k$; $gl>1$: puede utilizarse la prueba χ^2 si menos del 20% de las celdas registran valores menores que 5 en las frecuencias esperadas, y si no hay celdas con frecuencias esperadas menor que 1. Si estos requisitos no son cubiertos por los datos en la forma en que originalmente fueron recolectados, se recomienda combinar las categorías adyacentes para incrementar las frecuencias esperadas en las celdas. Esto ajusta los valores observados de χ^2 a la distribución teórica.

La prueba de la mediana

La distribución χ^2 ofrece la posibilidad de conocer si existen diferencias entre grupos utilizando la mediana de dos distribuciones de datos. Si bien existen diferentes pruebas estadísticas que sirven a este fin, un simple arreglo de una tabla de 2×2 permite utilizar la prueba χ^2 como una alternativa. Supongamos entonces que se tienen dos muestras independientes; entonces, se puede someter a prueba la hipótesis

de que ambas provienen de una misma población, o bien que ambas provienen de poblaciones diferentes.

Si se trata de demostrar que existen diferencias entre los grupos, la hipótesis de investigación establecerá que las medianas de ambas muestras son representativas de poblaciones diferentes. Por ende, la hipótesis nula establecerá que las medianas de las muestras serán las mismas (o estarán muy próximas entre sí), dado que las muestras provienen de la misma población. La hipótesis de investigación y la hipótesis nula se formalizan de la siguiente manera, donde los subíndices 1 y 2 representan las muestras respectivamente:

$$H_1: Mdn_1 \neq Mdn_2$$

$$H_0: Mdn_1 = Mdn_2$$

En este caso, la hipótesis de investigación es bidireccional, pero si el investigador sospecha que la diferencia observada se manifestará en alguna dirección predecible, puede formalizarse de esta manera:

$$H_1: Mdn_1 > Mdn_2$$

Aquí el investigador sostiene que la mediana del grupo 1 resultara significativamente mayor que la mediana del grupo 2.

La mediana combinada

La prueba de la mediana consiste en calcular la mediana combinada de las dos muestras ($Mdn_{1,2}$); esto se logra formando un solo grupo con todas las observaciones y calculando luego la mediana de ese grupo. Esa mediana se utiliza como punto de corte para separar los casos que caen por encima y por debajo de la mediana combinada, según prevengan del grupo 1 o del grupo 2. Este arreglo en una tabla 2x2 se representa de la siguiente manera:

	Grupo 1	Grupo 2
Valores por encima de la $Mdn_{1,2}$.	A	B
Valores por debajo de la $Mdn_{1,2}$	C	D

La celda A corresponde a los casos que están por encima de la mediana combinada y que pertenecen al grupo 1, la celda B corresponde a los casos que están por encima de la mediana combinada y que pertenecen al grupo 2. Lo mismo para las celdas C y D, que son aquellos casos que están por debajo de la mediana, según pertenezcan al

grupo 1 y 2. Si se cumple la situación anticipada en la hipótesis nula, la mediana de los grupos será la misma o estarán próximas entre sí. Por lo tanto, los valores por encima y por debajo de la mediana estarán proporcionalmente repartidos en las celdas A, B, C y D. Ahora bien, si la mediana de uno de los grupos es mayor que la mediana del otro, aquellos casos que muestren valores por encima de la mediana combinada tenderán a agruparse en la celda A o B. Bajo esta situación la distribución no será proporcional. En otras palabras, si hay diferencia entre los grupos se acumularán más casos en la celda correspondiente a la categoría valores por encima de la mediana combinada, y viceversa. Así se apreciará la concentración de casos opuestos por el vértice, que es propia cuando los atributos están asociados. Bajo tales circunstancias, el estadístico χ^2 será significativo. Por lo tanto, si dos muestras difieren en cuanto a su mediana, el procedimiento descrito permite utilizar una tabla de 2x2 y calcular un valor de χ^2 que si resulta significativo, se utiliza para rechazar la hipótesis nula que postula que las medianas son iguales por haber sido extraídas las muestras de la misma población.

Veamos ahora un ejemplo sencillo de la prueba de la mediana: en una escuela se toma un examen final de ciencias naturales a todos los alumnos de cuarto grado. Los maestros observan que en general, las notas de los alumnos del turno mañana son más altas que las del turno tarde y deciden averiguar si esta tendencia es producto del azar o existe alguna influencia que puede explicar las diferencias. Puesto que el examen no ha sido estandarizado, no se tiene certeza de que el promedio sea el estadístico adecuado para comparar a los alumnos, por ello deciden utilizar la prueba de la mediana. Los alumnos de cuarto grado en la escuela son 78 en total, y de acuerdo al arreglo de la tabla presentada anteriormente, se definió a los alumnos del turno mañana como Grupo 1, y a los alumnos del turno tarde como Grupo 2. La tabla quedó conformada como se muestra a continuación:

	Grupo 1	Grupo 2	
Valores por encima de la Mdn _{1,2} .	29	9	38
Valores por debajo de la Mdn _{1,2}	4	36	40
	33	45	78

Como se observa a simple vista, los valores de las celdas están opuestos por el vértice, lo que denota en principio que las medianas de los grupos no son las mismas. En otras palabras, existe una mayoría de casos en el grupo 1 que están por encima de la mediana combinada, y existe una mayoría de casos del grupo 2 por debajo de la mediana combinada. Con un arreglo de celdas como este, una alternativa para el cálculo del estadístico χ^2 , es utilizando la siguiente fórmula:

A	B	A+B
C	D	C+D
A+C	B+D	N

$$\chi^2 = \frac{N * (A * D - B * C)^2}{(A + B) * (C + D) * (A + C) * (B + D)}$$

Reemplazando los valores en la fórmula se tiene que:

$$\chi^2 = \frac{78 * (29 * 36 - 9 * 4)^2}{(38) * (40) * (33) * (45)} = 35.111$$

Consultando los valores tabulados de χ^2 se tiene que el valor observado tiene una probabilidad $p < 0.05$, lo cual nos resultaría suficiente para desestimar la hipótesis nula de que ambas muestras provienen de la misma población. Esto último, en el contexto del ejemplo que nos ocupa, se interpreta de la siguiente manera: si no existiera diferencia entre los alumnos de cuarto grado de ambos turnos, la proporción de casos por encima y por debajo de la media combinada debería haber sido proporcional, y por ende el valor del estadístico χ^2 debería haber estado próximo a 0. Dado que las distribuciones de frecuencias al interior de la tabla no es proporcional, el estadístico χ^2 se aparta notablemente de cero, superando su valor crítico con una probabilidad de ocurrencia, bajo la hipótesis nula, menor al 5%. Por ende es factible afirmar que los alumnos de cuarto grado del turno mañana muestran un rendimiento significativamente diferente a los alumnos de la tarde.

Análisis de una muestra simple: prueba χ^2 para la bondad de ajuste

La prueba χ^2 puede usarse para determinar si una distribución de frecuencias dadas se ajusta a los valores esperados de una distribución. En otras palabras, las diferencias entre las frecuencias observadas y las esperadas permiten cuantificar el grado en que una distribución empírica se aparta de una distribución teórica conocida. En muchos casos, la utilidad de la prueba χ^2 estriba en que permite verificar si el número de unidades de análisis que se distribuyen en las categorías de una variable, es verdaderamente diferente del esperado por azar. Es decir, es posible probar que existe una diferencia estadísticamente significativa entre el número observado de unidades de análisis que caen en cada categoría, y el número esperado que debería observarse

de ser cierta la hipótesis nula, que en este caso supone que las unidades de análisis se distribuyen aleatoriamente.

Al utilizar el estadístico χ^2 como una prueba de bondad de ajuste, las frecuencias esperadas necesarias para su cálculo se deducen del planteo de la hipótesis nula. Entonces, cuando se analiza una sola muestra distribuida en k categorías, las frecuencias esperadas se definen como:

$$f_e = N/k.$$

Esta prueba tiene ciertas restricciones en su uso, que dependen del valor de las frecuencias esperadas, así cuando $gl=1$ la frecuencia esperada debe ser igual o mayor que 5. Cuando $gl>1$, se debe comprobar que menos del 20% de las frecuencias esperadas sean menor que 5, y ninguna frecuencia esperada debe ser igual que 1.

Veamos un ejemplo sencillo: en una escuela se realizan evaluaciones periódicas de materiales didácticos con una serie de cinco listas de control de calidad. En un período de tiempo dado cada profesor debe tomar una de esas cinco listas al azar y realizar el control. Al final del mes se observó la cantidad de veces que se habían utilizado cada una de las listas, esperando una proporción cercana al azar en la cantidad de veces que se utilizaron. Teniendo en cuenta que durante ese mes se realizaron 45 controles de calidad, la frecuencia esperada para el uso de las cinco listas en 45 pruebas resulta en $45/5=9$. El número de veces que se utilizó cada lista se muestra en la tabla que sigue:

Listas	f observada	f esperada	Residual (fo-fe)
1	7	9	-2.0
2	16	9	7.0
3	6	9	-3.0
4	8	9	-1.0
5	8	9	-1.0
Total	45		

Teniendo en cuenta que la lista 2 se utilizó un número inesperadamente alto de veces, se sospechó que ella era la preferida por los profesores para realizar los controles sobre el material didáctico. Por ello se utilizó una prueba χ^2 para determinar si las frecuencias observadas podían ser representativas de una distribución por azar. El resultado se muestra en la siguiente tabla:

Variable Lista	
χ^2	7.111
gl	4
Sig (p)	0.130

El valor de significación crítico para χ^2 es de $p=0,05$. Puesto que el valor de significación observado supera al crítico, se infiere que el inusual aumento en el uso de la lista 2, está dentro de lo que corresponde esperar por azar. Nótese que en este caso no se realiza una comparación sobre valores observados y esperados de χ^2 puesto que la probabilidad de ocurrencia del estadístico es suficiente para determinar el ajuste a una distribución conocida.

Capítulo 9

Correlación de variables

Dos variables están relacionadas si los valores de una ellas se pueden predecir, en algún grado, de los valores observados en la otra. Esta idea puede comprenderse fácilmente si consideramos la altura y el peso de las personas; cuánto más alto es un individuo más pesado será. Esta variación conjunta entre la talla y el peso se puede estimar a partir de un coeficiente de correlación.

Determinar que existe correlación entre dos variables representa una importante fuente de conocimiento en materia de investigación, dado que podemos aprender mucho del comportamiento de una variable sobre la que sabemos muy poco, a partir del conocimiento que tenemos de otra variable de la que sabemos más. Comenzaremos con un ejemplo sencillo de este concepto. Supongamos que queremos averiguar qué factores o variables se relación con el promedio al final del primer año de la carrera. Para ello seleccionaremos una muestra aleatoria de estudiantes de una carrera universitaria (para este ejemplo el tamaño de la muestra será igual a 30), y recogemos datos sobre el promedio obtenido al final del primer año de estudio.

En principio suponemos que mientras más horas de estudio por semana le dedique el estudiante, más alto será su promedio. Por lo mismo suponemos que si el estudiante no puede disponer libremente del tiempo porque trabaja, estas horas de trabajo se restan a las horas de estudio, por lo tanto a mayor cantidad de horas trabajadas menor el rendimiento académico y más bajo el promedio. La historia académica del estudiante podría darnos alguna pista sobre el rendimiento, en tal caso tomamos el promedio obtenido al final de la secundaria como indicador de la trayectoria académica, suponiendo que aquellos estudiantes con mejor promedio al final del secundario, tendrán mejores promedios en la universidad. Finalmente, la ansiedad ante los exámenes es un factor que suele disminuir la competencia académica de estudiantes capaces, pero con dificultades para enfrentarlos. Por lo tanto, suponemos que aquellos alumnos con mayores niveles de ansiedad como rasgo de personalidad, serán los que muestren promedios más bajos.

Nótese que cada vez que suponemos la existencia de una relación entre variables también suponemos una dirección para esa relación. Por ejemplo, mayor cantidad de horas dedicadas al estudio más alto el promedio, implica que el aumento en la magnitud de una variable, es seguido por el aumento en la magnitud de la otra variable. Bajo estas circunstancias se dice que la relación entre las variables es directa. La relación también puede mostrar un sentido inverso, tal el caso en el que suponemos que altos niveles de ansiedad en el examen, dará como resultado un bajo rendimiento académico. Como se deduce, esto implica que el aumento en la magnitud

de una variable, es seguido por la disminución en la magnitud de la otra variable. Un coeficiente de correlación expresa mediante su signo esta cualidad de la relación entre variables. Así, un coeficiente de correlación positivo expresará una relación directa entre las variables, y un coeficiente negativo expresará una relación inversa.

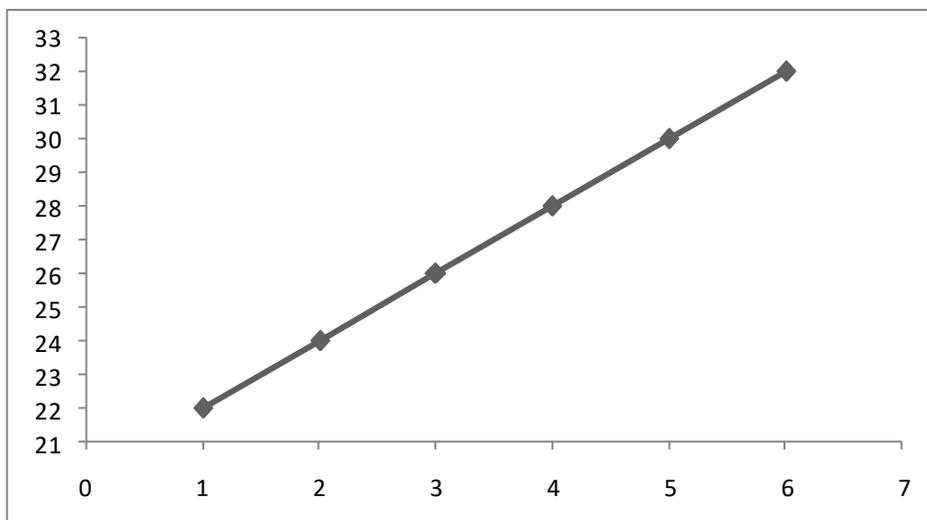
Asimismo, las variables se relacionan entre sí con una determinada magnitud o fuerza. Cuanto mayor sea la magnitud de la correlación, mayor determinación entre las variables. Por ejemplo, si el promedio al final del primer año dependiera exclusivamente de las horas de estudio, con solo saber cuántas horas le dedicó un individuo al estudio, sabríamos qué promedio obtuvo. Para determinar la magnitud y el tipo de relación conviene explorar mediante un gráfico denominado diagrama de dispersión.

El diagrama de dispersión

El diagrama de dispersión es un arreglo de las variables en ejes cartesianos que nos permite determinar qué tipo de relación existe entre ellas. En este caso nos ocuparemos de las relaciones lineales entre las variables. Comenzaremos analizando el tipo de relación entre dos variables ficticias x e y , definida por la relación $y = mx + k$, siendo m igual a 2 y k igual a 20.

X	Y
1	22
2	24
3	26
4	28
5	30

Como puede verse, para obtener un valor de y , todo lo que debemos hacer es sustituir en la ecuación el valor de x , por lo tanto, es sencillo construir la tabla para cualquier valor de x . La función $y = mx + k$, es la ecuación de una recta, y mediante ella todos los valores de y están perfectamente determinados por la función y por ende el diagrama cartesiano que pone en correspondencia x e y debería ser una recta, tal como se muestra en la siguiente figura.



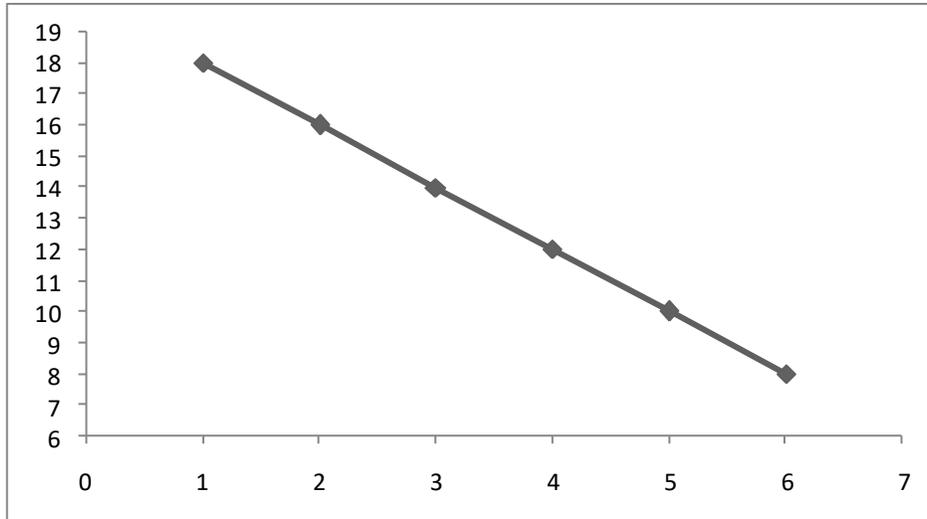
En el eje horizontal (abscisa) se han representado los valores de x y en el vertical (ordenadas) los valores de y . Como puede observarse la ecuación de la recta nos permite graficar justamente una recta. Lo importante para destacar aquí es que no estamos tratando con ninguna medición empírica, simplemente creamos una función de y cuyos valores dependen totalmente de x , por lo tanto, para saber cuál es el valor de y solo nos queda sustituir ese valor en la ecuación. Bajo tales circunstancias el cálculo de un coeficiente de correlación sería igual a 1, lo que equivale a decir que:

- a) la correlación entre las variables es perfecta, puesto que conociendo el valor de x , ya sabemos cuál es el valor de y ,
- b) puesto que el coeficiente es positivo, el incremento en el valor de x es seguido por un incremento en el valor de y .

Ahora vamos a utilizar los mismos valores de la tabla, pero la relación entre x e y estará dada por la función $y = -mx + k$, siendo m igual a -2 y k igual a 20 .

X	Y
1	18
2	16
3	14
4	12
5	10

La única diferencia con la función anterior es que ahora la constante m es negativa, por lo tanto la recta se oriente de manera inversa a la anterior, tal como lo muestra el siguiente gráfico:

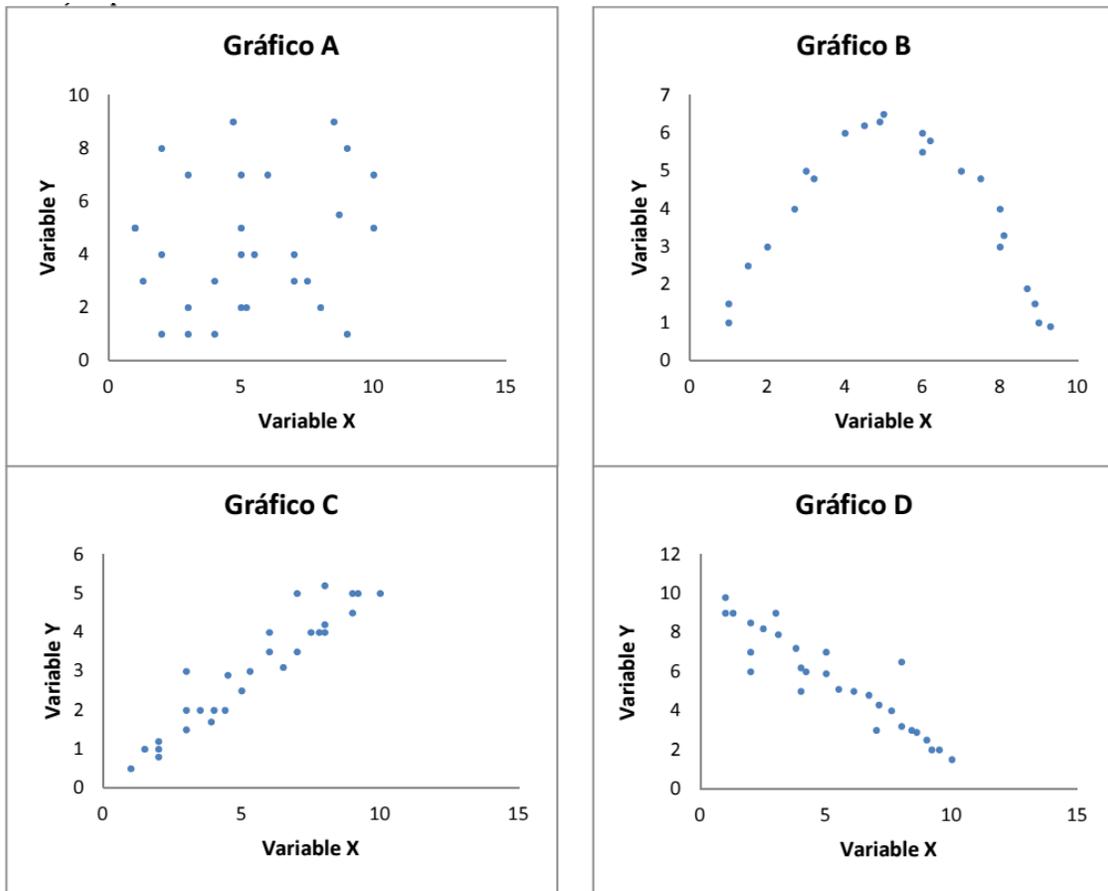


En este caso el cálculo de un coeficiente de correlación sería igual a -1 , lo que equivale a decir que:

- la correlación entre las variables es perfecta, puesto que conociendo el valor de x , ya sabemos cuál es el valor de y ,
- puesto que el coeficiente es negativo, el incremento en el valor de x es seguido por un decremento en el valor de y .

El diagrama de dispersión nos brinda una importante ayuda para determinar si dos variables empíricas están relacionadas de manera lineal. Si lo están, el coeficiente de correlación nos indica el sentido de la relación mediante el signo, y la fuerza de la relación mediante su magnitud.

Puede suceder que las variables no se relacionen de ninguna manera, en cuyo caso diremos que ambas son independientes. También puede ocurrir que la relación no sea lineal y deba expresarse mediante algún otro modelo matemático que no sea la ecuación de la recta. En la siguiente imagen se muestra un conjunto de situaciones posibles.



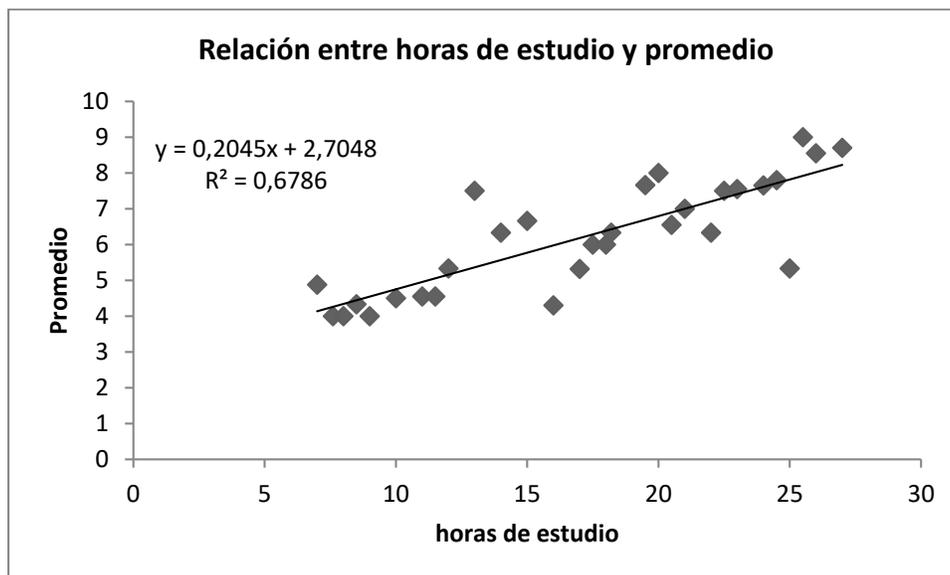
En la gráfica A se observa que los puntos cartesianos o pares ordenados x,y , están dispersos por todo el plano. Bajo estas circunstancias una predicción del comportamiento de la variable x a partir de los valores de la variable y es casi imposible, dicho en otras palabras, la predicción no es confiable. En la gráfica B se observa que existe un patrón en la distribución de los puntos cartesianos, pero dicho patrón no es descrito por una recta. El comportamiento de los pares ordenados describe una figura conocida como U invertida, que es propia de las funciones cuadráticas ($y=x^2$) que describen una parábola en su gráfica (dependiendo de la ecuación, la parábola puede abrirse hacia arriba en U, o hacia abajo así \cap). En ciencias sociales una de las curvas en U invertida más conocidas es la curva de Kuznets (basada en la hipótesis de la distribución propuesta por el economista Simon Kuznets), que puso en relación la distribución del ingreso (PBI) con la igualdad en su distribución.

Las gráficas C y D se asemejan a lo que ya vimos respecto de ecuaciones lineales, en donde C representa un modelo de relación positiva entre las variables y D un modelo de relación negativa. En ambos casos, la ecuación de la recta es la que mejor describe el comportamiento de las variables.

Coeficiente de correlación r de Pearson

Volvamos ahora a nuestro ejemplo, intentando determinar la relación entre las horas de estudio semanales y el rendimiento promedio. En primer lugar creamos un diagrama de dispersión para verificar el tipo de relación que existe entre las variables. En el ejemplo, se muestra un gráfico creado con Excel para las dos variables. Lo primero que se aprecia es que en el eje de las y se ha situado la variable promedio al final del primer año (promedio), en el eje de las x se ha situado la variable horas de estudio. Por cada valor de y existe un valor de x que ha sido representado en el plano cartesiano, dando origen a los puntos correspondientes a cada par $x;y$. Estos puntos forman lo que se denomina nube de puntos y nos da una primera impresión del tipo de relación entre las variables. Recuérdese que en nuestro ejemplo suponemos una relación directa entre las variables: a más horas de estudio mejor será el promedio. Como se aprecia en la disposición de los puntos en la nube, la relación estaría en consonancia con lo supuesto, ya que a más horas de estudio más alto es el promedio.

Si la relación fuera perfecta, todos los puntos se ubicarían sobre una recta, pero este raramente es el caso con las variables empíricas. En el gráfico se ha dibujado una recta teórica, que se ajusta a la nube de puntos de manera que todos ellos queden equidistante. Como se deduce, mientras más cercanos estén esos puntos a la línea recta, más próximo a la unidad estará el coeficiente de correlación.



En el ejemplo que nos ocupa el coeficiente de correlación r de Pearson, es $r=0.823$. Es un coeficiente positivo y próximo a uno, la relación entre las variables es directa y fuerte. Esto expresa que a mayor cantidad de horas de estudio semanales, mejor será

el promedio.

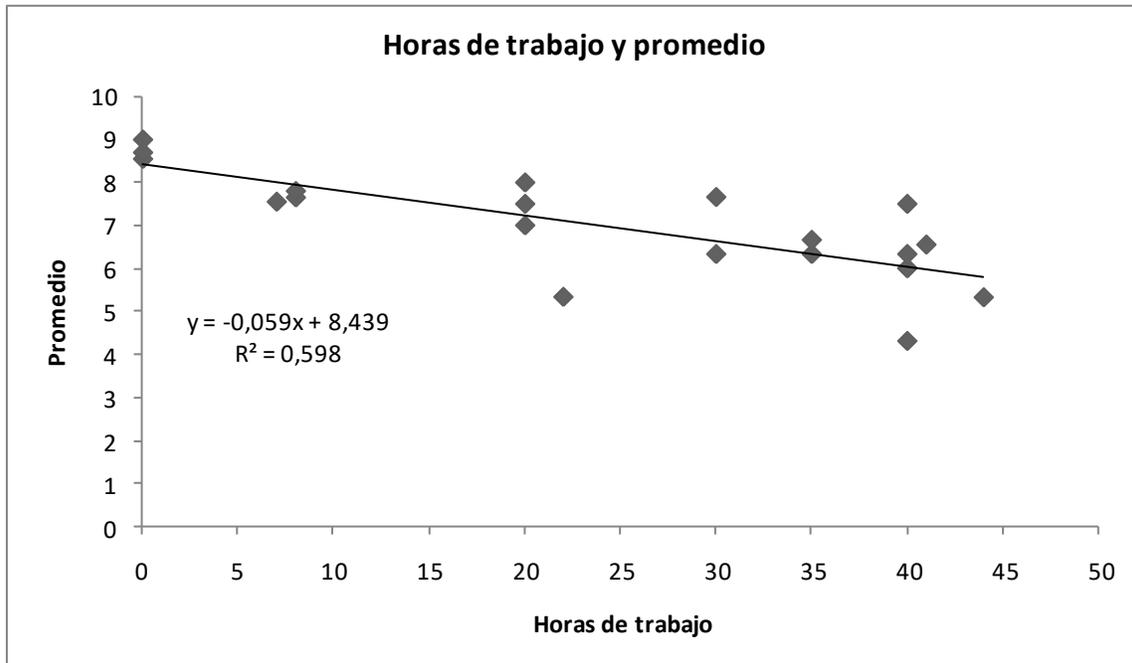
Puesto que la relación entre las variables es lineal, se puede expresar mediante la ecuación de la recta que en este caso es $y=0.204x+2.704$. Esta ecuación, denominada ecuación de regresión, permite calcular cualquier valor de x , pero debemos tener en cuenta que los valores de x calculados estarán sujetos a un cierto error de estimación, dado que la relación entre las variables no es igual a 1.

En el gráfico también se muestra otro coeficiente, llamado de determinación y que expresa el monto de varianza de y que puede explicarse a partir de x . Este coeficiente R^2 puede interpretarse a partir de su valor porcentual, que en nuestro ejemplo es de 67,8%.

Ahora sí estamos en condiciones de afirmar que aquellos alumnos que dedican más horas de estudio semanales tendrán un mejor promedio al final del año. Si nos preguntamos cuánto es lo que podemos explicar de la variación en los promedios de todos los individuos de nuestra muestra, utilizando las horas de estudio como variable independiente, podemos afirmar que aproximadamente el 68%. Finalmente, si nos preguntamos cuál sería aproximadamente el promedio de un individuo que estudiara 35 horas semanales, aplicamos la ecuación de regresión:

$$y= 0.204x35+2.704= 9,84$$

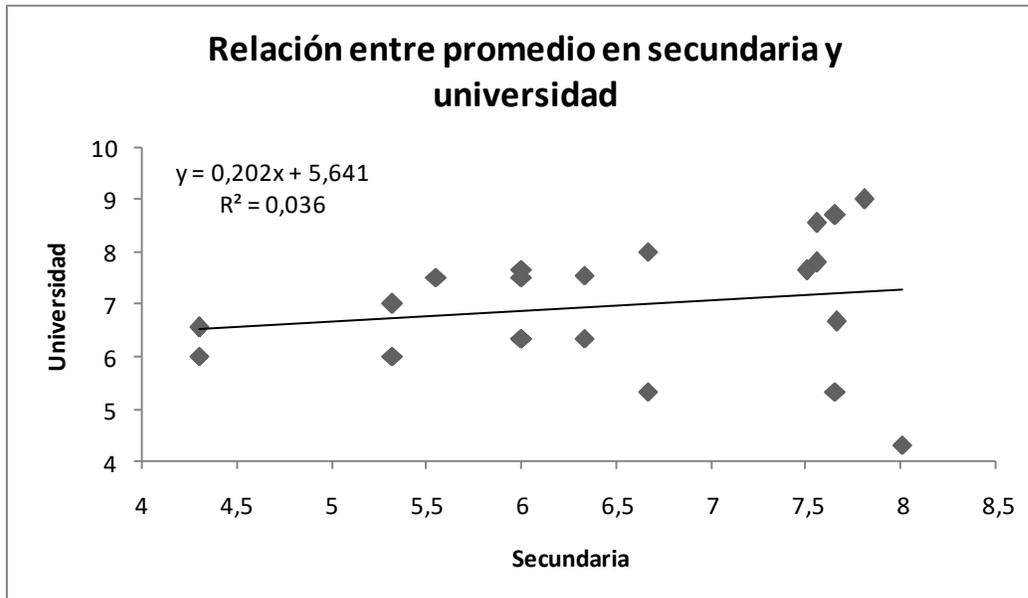
Otra suposición que hicimos al comienzo fue que si el estudiante no puede disponer libremente del tiempo porque trabaja, las horas de trabajo se restan a las horas de estudio, por lo tanto a mayor cantidad de horas trabajadas menor el rendimiento académico y más bajo el promedio. En la muestra de estudiantes encontramos algunos que no trabajan, otros que trabajan algunas horas semanales, y otros que lo hacen en jornadas de ocho o más horas semanales. En el diagrama de dispersión queda expresada la relación entre las variables y nos permite determinar cómo afecta el trabajo al rendimiento académico. Nuevamente, la relación se ha expresado mediante un gráfico de Excel.



El coeficiente de correlación de Pearson en este caso es $r = -0.773$. Es un coeficiente alto pero de signo negativo, lo cual implica una relación inversa entre las variables. El coeficiente pone en evidencia que a mayor cantidad de horas de trabajo, menor el promedio del estudiante. La cantidad de varianza que puede ser explicada del promedio utilizando las horas de trabajo como variable independientes es aproximadamente el 60%. Finalmente, si quisiéramos estimar cual sería el promedio de una persona que trabaja 25 horas semanales, utilizamos la ecuación de regresión:

$$y = -0.059 \times 25 + 8.439 = 6.964 \approx 7$$

Otra de las variables que nos proponíamos analizar, era la historia académica del estudiante suponiendo que aquellos estudiantes con mejor promedio al final del secundario, tendrán mejores promedios en la universidad. Procediendo de la misma manera en que lo hicimos en los casos anteriores, ponemos las variables en correspondencia en un diagrama de dispersión, tal como se muestra a continuación.



El coeficiente de correlación en este caso es $r=0.189$, un coeficiente positivo de baja magnitud. Esto se interpreta como una relación muy débil entre las variables. El promedio al final del secundario nos informa muy poco del rendimiento académico de un estudiante en universidad, aun mas, si quisiéramos explicar el rendimiento académico universitario del conjunto de datos a partir de la variable promedio en el secundario, solo alcanzaríamos a dar cuenta del 3,6% de la varianza.

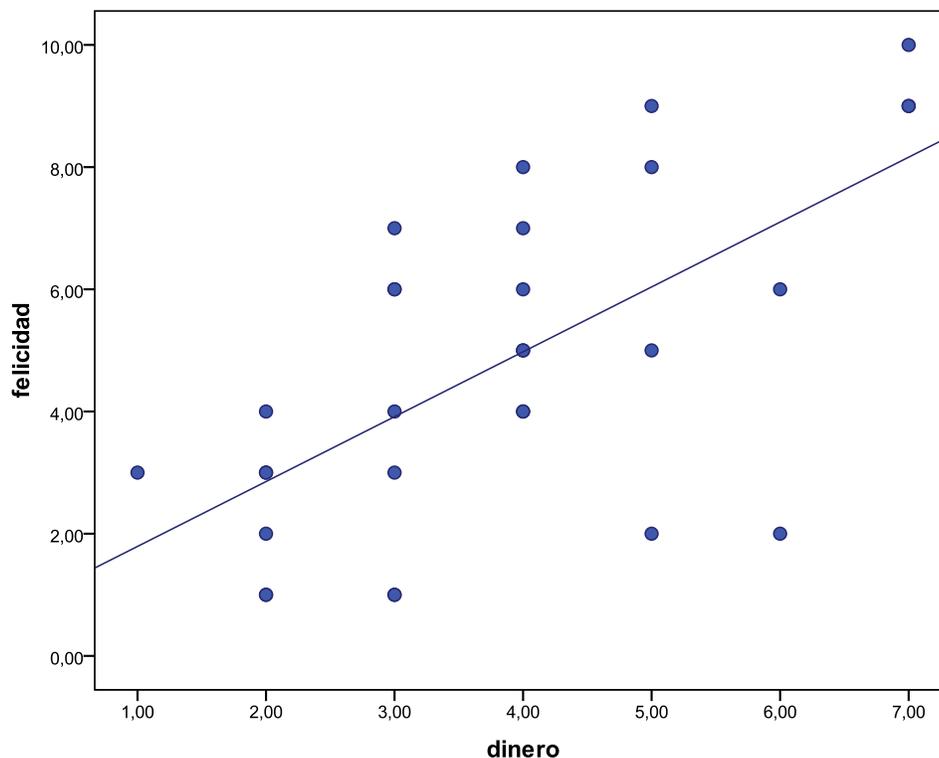
En este punto analizaremos más detenidamente el gráfico. Recordemos que en el caso de una asociación perfecta entre las variables, todos los puntos de la nube caen sobre la recta de regresión. A medida que el coeficiente se aleja de la unidad y se aproxima a cero la nube de puntos se hace cada vez más dispersa en torno a la recta. Esto significa que es cada vez más imprecisa alguna predicción de los valores de una variable a partir de los valores de la otra variable. Veamos ahora el gráfico, y se notará que los puntos se dispersan cada vez más a medida que ascendemos en los valores de x , que en este caso es promedio en secundaria. Si tomamos el valor promedio 4,33 en secundaria, notaremos que hay dos valores próximos entre sí en la variable promedio en la universidad, estos son 6 y 6,55, pero si tomamos el valor 7,55 en la variable promedio en secundaria notaremos que hay valores de 5 hasta 8,5 aproximadamente. En tal caso, la predicción que pudiéramos hacer sería demasiado ambigua e inexacta, dada la dispersión de los valores. Esta es la razón por la cual el coeficiente de correlación entre las variables tiene una magnitud tan baja.

Coeficiente de correlación r_s de Spearman

Charles Spearman durante sus investigaciones sobre la inteligencia humana, desarrollo un coeficiente de correlación basado en rangos, que hoy lleva su nombre. Este coeficiente de correlación es una opción cuando una o ambas variables se han medido en escala ordinal. Es decir, es una medida que puede aplicarse a aquellos casos en que las escalas no son métricas. Este coeficiente resulta de mucha utilidad cuando se trabajan con escalas de tipo Lickert, donde los individuos juzgan un atributo en una escala ordinal especialmente diseñada para el mismo.

Veamos un ejemplo: En un estudio se pidió a un grupo de personas que juzgaran simultáneamente su estado de ánimo actual y la importancia que ellas le otorgaban al dinero. El objetivo de la investigación consistió en determinar si existía relación entre el estado de ánimo y factores materiales. Para el estado de ánimo se utilizó una escala autoadministrada de felicidad en donde el individuo debía indicar en un orden de 0 a 10 cómo se sentía en ese momento, siendo 10 el nivel más alto y correspondiente a muy feliz. Igualmente para la importancia del dinero, se utilizó un orden de 0 a 7 donde el máximo puntaje representaba muy importante. Dado que ambas variables están diseñadas sobre una escala ordinal y no poseen la misma cantidad de reactivos, se utilizó el coeficiente de correlación r_s para determinar el grado de correlación entre las variables.

Como en los casos anteriores analizaremos primero el diagrama de dispersión. En este caso ha sido graficado sobre un total de 30 casos para cada variable y utilizando el programa SPSS.



Como puede observarse, el diagrama sugiere una correlación lineal positiva, de manera que aquellos individuos que juzgan el dinero como más importante, son los que tienden a mostrarse más felices. Nótese sin embargo que existen personas con altos niveles de felicidad, que no le dan al dinero demasiada importancia, juzgándolo en valores bajos. Estos individuos se apartan de la tendencia general. Además, se evidencia que hay individuos que le otorga un valor de importancia muy alto al dinero, pero no manifiesta un alto grado de felicidad. Con estos indicios es de suponer que el coeficiente de correlación tendrá un valor medio-alto. El cálculo del coeficiente nos muestra que efectivamente este alcanza el valor $r_s=0.619$.

El coeficiente de correlación ordinal de Spearman nos permite trabajar con variables ordinales y su interpretación es la misma que ya vimos para el coeficiente de correlación de Pearson. Sin embargo no es conveniente con este tipo de coeficiente calcular una ecuación de regresión o el coeficiente de determinación R^2 , dada la naturaleza ordinal de las variables.

Coeficiente de correlación y prueba de hipótesis

Se puede plantear como hipótesis que dos variables están relacionadas, en tal caso, el planteo implica determinar la magnitud y el signo de la correlación encontrada entre las variables con una distribución teórica que estipula que ambas variables no tienen

relación. Esto último corresponde a la hipótesis nula. La prueba de hipótesis se basa en un resultado sobre el cual rechazaremos o no ésta hipótesis, basándonos en la probabilidad de encontrar un resultado específico, considerando cierta la hipótesis nula. Un ejemplo nos permitirá clarificar los conceptos expuestos. Basándose en estudios previos, un investigador supone que la madurez escolar está relacionada con las habilidades lingüísticas. Para probar esta afirmación evalúa una muestra de 12 escolares con dos pruebas específicas, una para cada variable. Las pruebas utilizadas están diseñadas de modo que sus puntuaciones responden a una escala métrica, por lo cual el investigador decide comprobar la relación utilizando el coeficiente de correlación de Pearson. Los datos recogidos se muestran en la siguiente tabla:

Sujetos	Madurez Escolar (x)	Habilidades lingüísticas (y)
1	10	15
2	14	19
3	9	14
4	8	13
5	13	18
6	13	19
7	16	18
8	6	9
9	11	21
10	7	16
11	5	12
12	15	20

La tabla muestra en la primera columna los sujetos que participaron en el estudio, la segunda columna muestra los puntajes obtenidos en madurez escolar que hemos denominado como x , la tercera columna muestra los puntajes en la prueba de habilidades lingüísticas que hemos denominado y . La hipótesis nula implica negar la hipótesis de investigación, por lo cual su planteo implica que no existe relación entre las variables. Ahora corresponde expresar en términos estadísticos ambas hipótesis:

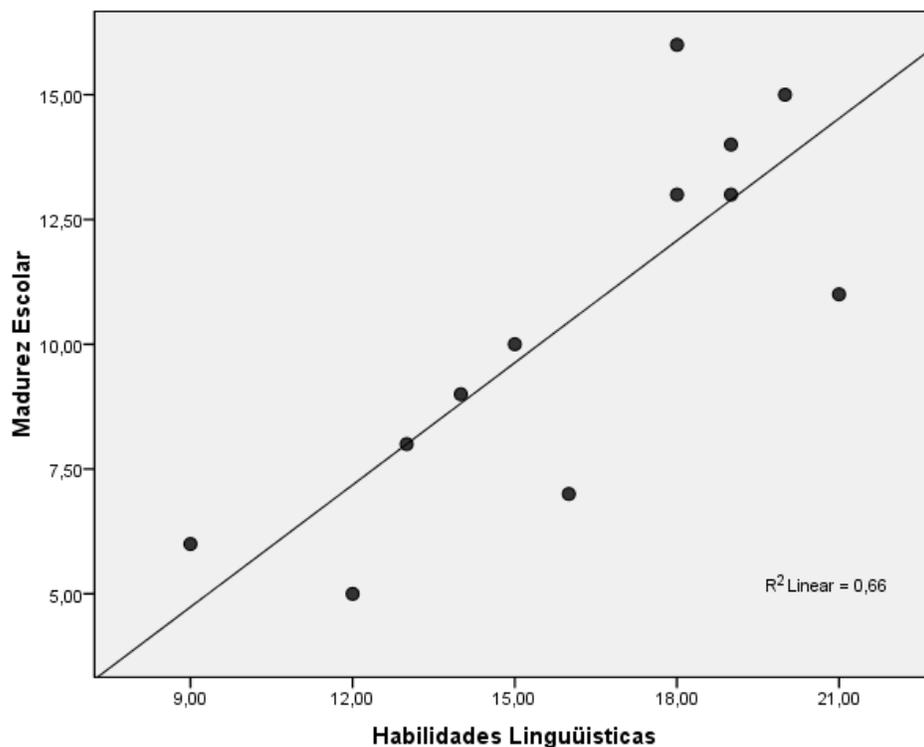
$H_0: r=0$

$H_1: r \neq 0$

Nótese que la formulación estadística supone que el coeficiente de correlación será

cero bajo el supuesto de la hipótesis nula. Al contrario, la hipótesis de investigación supone un coeficiente de correlación diferente de cero. Corresponde ahora determinar bajo qué circunstancias es factible rechazar la hipótesis nula. Pues bien, suponiendo que en la población el coeficiente de correlación para las variables mencionadas es cero, la probabilidad de encontrar un coeficiente de una magnitud distinta de cero resulta muy poco probable. Bajo estas condiciones establecemos el valor de α correspondiente a la probabilidad de cometer un error tipo I. Usualmente este valor es $\alpha=0.05$.

El diagrama de dispersión y el cálculo del valor del coeficiente de correlación, junto a su probabilidad de ocurrencia bajo la hipótesis nula, se realizó utilizando el programa SPSS.



El cálculo del coeficiente de correlación de Pearson $r=0,812$, y la probabilidad asociada a este coeficiente bajo la hipótesis nula es $p=0,001$. Ahora bien, considerando que la probabilidad de haber obtenido un coeficiente de correlación como el reportado es muy baja tomando en cuenta la hipótesis nula, resulta que el supuesto de que ambas variables no tienen relación debe ser rechazado. En otras palabras, rechazamos la hipótesis nula que supone que no existe relación entre las variables, a favor de la hipótesis que sostiene que las variables están relacionadas.

Finalmente vemos que el coeficiente de determinación, tomando como variable

independiente a las habilidades lingüísticas, es de 0,66. Como ya dijimos, esto equivale a que el 66% de la varianza en la madurez escolar puede explicarse mediante las habilidades lingüísticas del individuo.

Veamos otro ejemplo de prueba de hipótesis, donde la hipótesis de investigación es más específica: bajo el supuesto que la memoria es un recurso directamente involucrado en el desarrollo de las capacidades verbales, en un experimento se analizó si la amplitud de la memoria verbal inmediata, se relacionaba con la capacidad de resolver el cierre de sentido de oraciones sintácticamente ambiguas. Los investigadores utilizaron la prueba de retención de dígitos para medir la amplitud de memoria verbal, la cual consiste en repetir de manera correcta series de números no consecutivos que van 2 a 9 dígitos. Por su parte, la capacidad de resolución de ambigüedades se estudió mediante la aplicación de la prueba Cloze, en la cual se asigna un puntaje de cero a veinte, de acuerdo a la cantidad de respuestas correctas. En la etapa inicial del estudio, se tomó una muestra de diez estudiantes universitarios y los resultados de la aplicación de estas pruebas se detallan en la siguiente tabla. Si bien la prueba de Cloze es una prueba estándar, sus puntajes no se consideran de escala métrica, sino ordinal. Por lo tanto, el coeficiente de correlación que se emplea para la prueba de hipótesis es r_s .

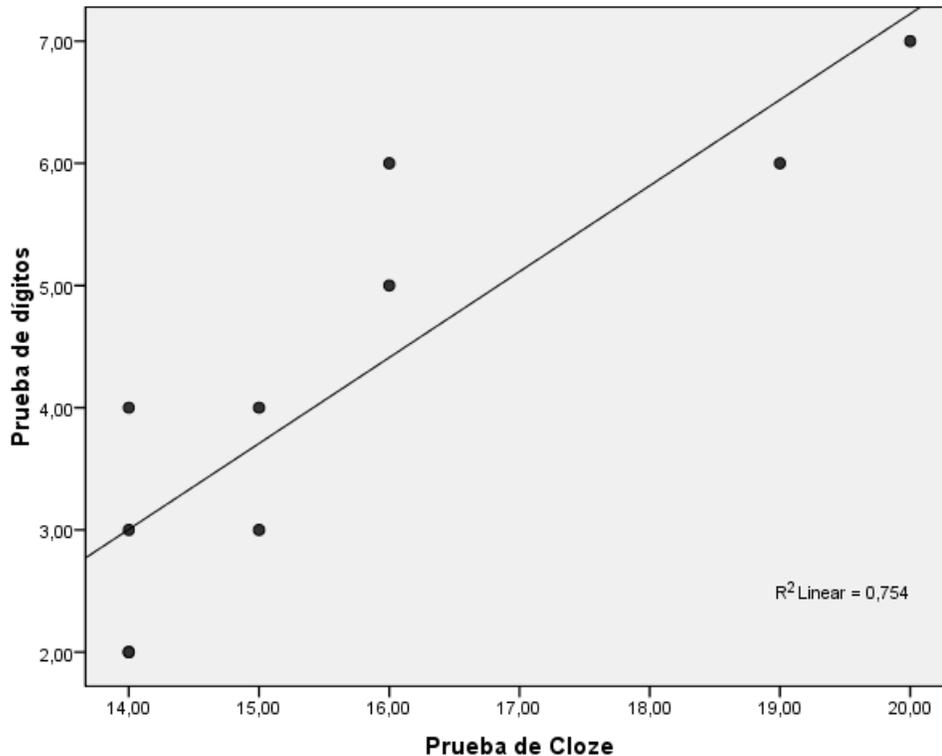
casos	Prueba de dígitos (x)	Prueba de Cloze (y)
1	2	17
2	3	17
3	2	14
4	4	18
5	3	13
6	4	14
7	5	15
8	6	16
9	6	15
10	7	15

Dado que los investigadores suponen que la memoria está directamente involucrada como recurso cognitivo en la prueba de resolución de ambigüedades, se espera encontrar una correlación directa (positiva) entre las variables. La formulación estadística de la hipótesis de investigación y la hipótesis nula sería:

$H_0: r=0$

$H_1: r>0$

Seguidamente analizamos el diagrama de dispersión y se calcula el valor del coeficiente de correlación, junto a su probabilidad de ocurrencia bajo la hipótesis nula. En este ejemplo también, utilizamos el programa SPSS para realizar esos cálculos.



El cálculo del coeficiente de correlación de Spearman $r_s=0,889$, y la probabilidad asociada a este coeficiente bajo la hipótesis nula es $p<0,001$. Aquí también resulta que la probabilidad de obtener un coeficiente de correlación como el reportado es muy baja, por lo tanto se rechaza la hipótesis que ambas variables no tienen relación. En otras palabras, rechazamos la hipótesis nula. El coeficiente de determinación, tomando como variable independiente a la amplitud de memoria verbal, es de 0,754, lo cual representa una explicación de la varianza en la prueba de Cloze del 75,4%.

Algunas precisiones sobre el coeficiente de correlación

Un coeficiente de correlación indica la magnitud de una relación lineal, a mayor valor de r , mayor grado de correlación, por ende si se compara una correlación de $r=0.2$ con una $r=0.4$, se concluye que esta última es mayor que la primera, pero no es correcto afirmar que es el doble. Si se desea comparar dos o más coeficientes se utiliza el cuadrado de la correlación r , esto es R^2 que como ya vimos, representa la proporción de la varianza explicada en la variable dependiente; esto se denomina reducción proporcional del error.

De acuerdo a lo dicho, para comparar dos coeficientes de correlación como $r_1=0.2$, y $r_2=0.4$, las elevamos al cuadrado, y entonces tenemos que $R^2_1= 0.04$ y $R^2_2=0.16$. En proporción, la magnitud de la segunda correlación es cuatro veces mayor que la primera y no dos veces como se deduciría si usáramos el coeficiente de correlación r .

¹ <http://www.cronista.com/economiapolitica/Recorre-el-mapa-interactivo-del-delito-en-Argentina-el-Gobierno-publico-estadisticas-comparativas-20160425-0089.html>