

Estadística

Alejandro Cáceres (alejandro.caceres.dominguez@upc.edu)

2023-02-16

Contents

1	Objetivo	5
1.1	Lectura recomendada	6
2	Descripción de datos	9
2.1	Método científico	9
2.2	Estadística	9
2.3	Datos	10
2.4	Tipos de resultado	10
2.5	Experimentos aleatorios	10
2.6	Frecuencias absolutas	11
2.7	Frecuencias relativas	11
2.8	Diagrama de barras	12
2.9	Gráfico de sectores (pie)	13
2.10	Variables categóricas ordinales	13
2.11	Frecuencias acumuladas absolutas y relativas	14
2.12	Gráfica de frecuencia acumulada	15
2.13	Variables numéricas	15
2.14	Transformando datos continuos	16
2.15	Tabla de frecuencias para una variable continua	17
2.16	Histograma	17
2.17	Gráfica de frecuencia acumulada	19
2.18	Estadísticas de resumen	19
2.19	Promedio (media muestral)	20
2.20	Promedio	21
2.21	mediana	21
2.22	Dispersión	23
2.23	Variación de la muestra	24
2.24	Rango intercuartílico (IQR)	25
2.25	Diagrama de caja	26
2.26	Preguntas	27
2.27	Ejercicios	28
3	Probabilidad	29

3.1	Experimentos aleatorios	29
3.2	Probabilidad de medición	30
3.3	Probabilidad clásica	30
3.4	Frecuencias relativas	30
3.5	Frecuencias relativas en el infinito	32
3.6	Probabilidad frecuentista	33
3.7	Probabilidades clásicas y frecuentistas	34
3.8	Definición de probabilidad	34
3.9	Tabla de probabilidades	35
3.10	Espacio muestral	35
3.11	Eventos	36
3.12	Álgebra de eventos	36
3.13	Resultados mutuamente excluyentes	36
3.14	Probabilidades conjuntas	37
3.15	Tabla de contingencia	38
3.16	La regla de la suma:	38
3.17	Preguntas	39
3.18	Ejercicios	40

Chapter 1

Objetivo

Este es el curso de introducción a la estadística de la EEBE (UPC).

La estadística es un **lenguaje** que permite afrontar problemas nuevos, sobre los que no tenemos solución, y en donde interviene la **aleatoriedad**.

En este curso trataremos los **conceptos fundamentales** de estadística.

- 3 horas de **teoría** por semana: Explicaremos los conceptos, haremos ejercicios.
- 6 horas de **estudio individual** por semana: Notas de curso y los recursos en ATENEA.
- 2 horas de Solución de problemas con **R**: Sesiones presenciales con ordenador (Prácticas).

Las fechas de exámenes y material de estudio adicional se pueden encontrar en **ATENEA metacurso**:

Objetivos de evaluación:

Q1 (10%): Prueba en ordenador duración 2h en las fechas indicadas.

- a. Dominio de comandos básicos en R (Prácticas)
- b. Capacidad de calcular estadísticos descriptivos y gráficos, en situaciones concretas (Teoría/Práctica)
- c. Conocimiento sobre la regresión lineal (Prácticas)

EP1 (25%): Prueba escrita (2-3 problemas)

- a. Capacidad de interpretación de enunciados en fórmulas de probabilidad (Teoría).
- b. Conocimiento de las herramientas básicas para solucionar problemas de probabilidad conjunta y probabilidad condicional (Teoría).

- c. Dominio matemático de funciones de probabilidad para calcular sus propiedades básicas (Teoría).

Q2 (10%): Prueba en ordenador duración 2h en las fechas indicadas

- a. Capacidad de identificación de modelos de probabilidad en problemas concretos (Teoría/Práctica).
- b. Uso de funciones de R para calcular probabilidades de modelos probabilísticos (Práctica/Teoría)

Q3 (10%): Prueba en ordenador duración 2h en las fechas indicadas

- a. Capacidad de identificación de un estadístico de muestreo y sus propiedades (Teoría/Práctica)
- b. Conocimiento de cómo calcular la probabilidad de los estadísticos de muestreo (Teoría/Práctica)
- c. Uso de comandos en R para calcular probabilidades y hacer simulaciones de muestras aleatorias (Prácticas)

EP2 (10%): Prueba escrita (2-3 problemas)

- a. Capacidad matemática para determinar estimadores puntuales de modelos de probabilidad.
- b. Conocimiento de las propiedades de los estimadores puntuales.

CG (5%): Prueba escrita (2 preguntas sobre un texto)

- a. Capacidad de expresión escrita sobre un tema relacionado a la estadística.

EP3 (30%): Prueba por ordenador presencial (2-3 problemas)

- a. Conocimiento de los intervalos de confianza y sus propiedades (Teoría).
- b. Capacidad de identificar el tipo de intervalo de confianza en un problema concreto (Teoría).
- c. Capacidad de interpretación del tipo de hipótesis a usar en un problema concreto (Teoría).
- d. Propiedades de las pruebas de hipótesis.
- e. Uso de comandos en R para resolver problemas de intervalos de confianza y pruebas de hipótesis (Práctica).

coordinadores:

- Luis Mujica (luis.eduardo.mujica@upc.edu)
- Pablo Buenestado (pablo.buenestado@upc.edu)

1.1 Lectura recomendada

- Las notas de clase de nuestra sección estarán accesibles en ATENEA en pdf y en html.

- Douglas C. Montgomery and George C. Runger. “Applied Statistics and Probability for Engineers” 4th Edition. Wiley 2007.

Chapter 2

Descripción de datos

En este capítulo, presentaremos herramientas para describir datos.

Lo haremos utilizando tablas, figuras y estadísticos descriptivos de tendencia central y dispersión.

También presentaremos conceptos clave en estadística como experimentos aleatorios, observaciones, resultados y frecuencias absolutas y relativas.

2.1 Método científico

Uno de los objetivos del método científico es proporcionar un marco para resolver los problemas que surgen en el estudio de los fenómenos naturales o en el diseño de nuevas tecnologías.

Los humanos modernos han desarrollado un **método** durante miles de años que todavía está en desarrollo.

El método tiene tres actividades humanas principales:

- *Observación* caracterizada por la adquisición de **datos**
- *Razón* caracterizada por el desarrollo de **modelos** matemáticos
- *Acción* caracterizada por el desarrollo de nuevos **experimentos** (tecnología)

Su compleja interacción y resultados son la base de la *actividad científica*.

2.2 Estadística

La estadística se ocupa de la interacción entre *modelos* y *datos* (la parte inferior de la figura).

Las preguntas de tipo estadístico son:

- ¿Cuál es el mejor modelo para mis datos (inferencia)?
- ¿Cuáles son los datos que produciría un determinado modelo (predicción)?

2.3 Datos

Los datos se presentan en forma de observaciones.

Una **Observación** o *Realización* es la adquisición de un número o una característica de un experimento.

Por ejemplo, tomemos la serie de números que se producen por la repetición de un experimento (1: éxito, 0: fracaso)

... 1 0 0 1 0 1 0 1 1 ...

El número en negrita es **una observación** en una repetición del experimento

Un **resultado** es una **posible** observación que es el resultado de un experimento.

1 es un resultado, **0** es el otro resultado del experimento.

Recuerda que la observación es **concreta** es el número que obtienes un día en el laboratorio. El resultado **abstracto** es una de las características del tipo de experimento que estás realizando.

2.4 Tipos de resultado

En estadística nos interesan principalmente dos tipos de resultados.

- **Categoricos:** Si el resultado de un experimento es una cualidad. Pueden ser nominales (binario: sí, no; múltiple: colores) u ordinales cuando las cualidades pueden jerarquizarse (gravedad de una enfermedad).
- **Numéricos:** Si el resultado de un experimento es un número. El número puede ser discreto (número de correos electrónicos recibidos en una hora, número de leucocitos en sangre) o continuo (estado de carga de la batería, temperatura del motor).

2.5 Experimentos aleatorios

Se puede decir que el tema de estudio de la estadística son los experimentos aleatorios, el medio por el cual producimos datos.

Definición:

Un **experimento aleatorio** es un experimento que da diferentes resultados cuando se repite de la misma manera.

Los experimentos aleatorios son de diferentes tipos, dependiendo de cómo se realicen:

- en el mismo objeto (persona): temperatura, niveles de azúcar.
- sobre objetos diferentes pero de la misma medida: el peso de un animal.
- sobre eventos: el número de huracanes por año.

2.6 Frecuencias absolutas

Cuando repetimos un experimento aleatorio con resultados **categoricos**, registramos una lista de resultados.

Resumimos las observaciones contando cuántas veces vimos un resultado particular.

Frecuencia absoluta:

$$n_i$$

es el número de veces que observamos el resultado i .

Ejemplo (leucocitos)

Extraigamos un leucocito de **un** donante y anotemos su tipo. Repitamos el experimento $N = 119$ veces.

(célula T, célula T, neutrófilo, ..., célula B)

La segunda **célula T** en negrita es la segunda observación. La última **célula B** es la observación número 119.

Podemos listar los **resultados** (categorías) en una **tabla de frecuencia**:

```
##      outcome ni
## 1      T Cell 34
## 2      B cell 50
## 3  basophil 20
## 4  Monocyte  5
## 5 Neutrophil 10
```

De la tabla, podemos decir que, por ejemplo, $n_1 = 34$ es el número total de células T observadas en la repetición del experimento. También notamos que el número total de repeticiones $N = \sum_i n_i = 119$.

2.7 Frecuencias relativas

También podemos resumir las observaciones calculando la **proporción** de cuántas veces vimos un resultado en particular.

$$f_i = n_i/N$$

donde N es el número total de observaciones

En nuestro ejemplo se registraron $n_1 = 34$ células T, por lo que nos preguntamos por la proporción de células T del total de 119. Podemos agregar estas proporciones f_i en la tabla las frecuencias.

```
##      outcome ni      fi
## 1      T Cell 34 0.28571429
## 2      B cell 50 0.42016807
## 3    basophil 20 0.16806723
## 4    Monocyte  5 0.04201681
## 5 Neutrophil 10 0.08403361
```

Las frecuencias relativas son **fundamentales** en estadística. Dan la proporción de un resultado en relación con los otros resultados. Más adelante las entenderemos como las observaciones de las probabilidades.

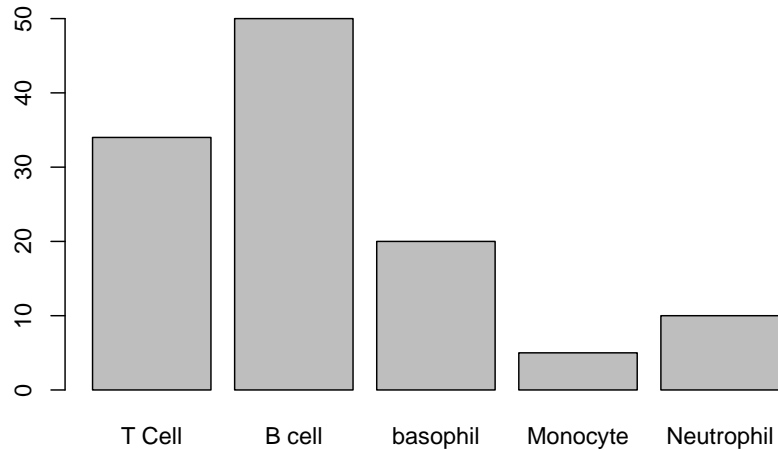
Para las frecuencias absolutas y relativas tenemos las propiedades

- $\sum_{i=1..M} n_i = N$
- $\sum_{i=1..M} f_i = 1$

donde M es el número de resultados.

2.8 Diagrama de barras

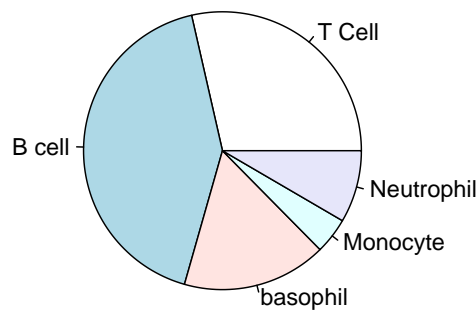
Cuando tenemos muchos resultados y queremos ver cuáles son los más probables, podemos usar un gráfico de barras que es una cifra de n_i Vs los resultados.



2.9 Gráfico de sectores (pie)

También podemos visualizar las frecuencias relativas con un gráfico de sectores.

El área del círculo representa el 100% de las observaciones (proporción = 1) y las secciones las frecuencias relativas de cada resultado.



2.10 Variables categóricas ordinales

El tipo de leucocito de los ejemplos anteriores es una variable nominal **categórica**. Cada observación pertenece a una categoría (cualidad). Las categorías no siempre tienen un orden determinado.

A veces, las variables **categóricas** se pueden **ordenar** cuando cumplen una clasificación natural. Esto permite introducir **frecuencias acumulativas**.

Ejemplo (misofonía)

Este es un estudio clínico en 123 pacientes que fueron examinados por su grado de misofonía. La misofonía es ansiedad/ira descontrolada producida por ciertos sonidos.

Cada paciente fue evaluado con un cuestionario (AMISO) y se clasificaron en 4 grupos diferentes según la gravedad.

Los resultados del estudio son

```
##      [1] 4 2 0 3 0 0 2 3 0 3 0 2 2 0 2 0 0 3 3 0 3 3 2 0 0 0 4 2 2 0 2 0 0 0 3 0 2
```

```
## [38] 3 2 2 0 2 3 0 0 2 2 3 3 0 0 4 3 3 2 0 2 0 0 0 2 2 0 0 2 3 0 1 3 2 4 3 2 3
## [75] 0 2 3 2 4 1 2 0 2 0 2 0 2 2 4 3 0 3 0 0 0 2 2 1 3 0 0 3 2 1 3 0 4 4 2 3 3
## [112] 3 0 3 2 1 2 3 3 4 2 3 2
```

Cada observación es el resultado de un experimento aleatorio: medición del nivel de misofonía en un paciente. Esta serie de datos se puede resumir en términos de los resultados en la tabla de frecuencia

```
## outcome ni      fi
## 1      0 41 0.33333333
## 2      1  5 0.04065041
## 3      2 37 0.30081301
## 4      3 31 0.25203252
## 5      4  9 0.07317073
```

2.11 Frecuencias acumuladas absolutas y relativas

La gravedad de la misofonía es **categorica ordinal** porque sus resultados pueden ordenarse en relación con su grado.

Cuando los resultados se pueden ordenar, es útil preguntar cuántas observaciones se obtuvieron hasta un resultado dado. Llamamos a este número la **frecuencia acumulada absoluta** hasta el resultado i :

$$N_i = \sum_{k=1..i} n_k$$

También es útil para calcular la **proporción** de las observaciones que se obtuvo hasta un resultado dado

$$F_i = \sum_{k=1..i} f_k$$

Podemos agregar estas frecuencias en la **tabla de frecuencias**

```
## outcome ni      fi  Ni      Fi
## 0      0 41 0.33333333 41 0.33333333
## 1      1  5 0.04065041 46 0.3739837
## 2      2 37 0.30081301 83 0.6747967
## 3      3 31 0.25203252 114 0.9268293
## 4      4  9 0.07317073 123 1.0000000
```

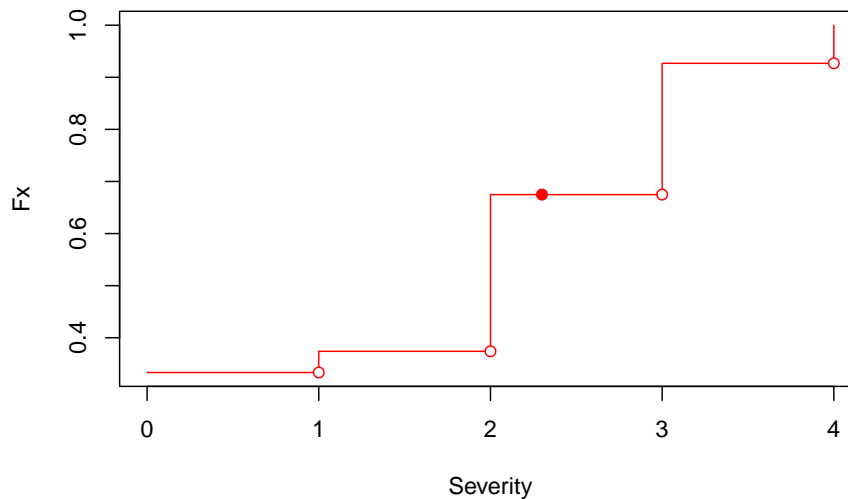
Por lo tanto, el **67 %** de los pacientes tenían misofonía hasta la gravedad **2** y el **37 %** de los pacientes tenían una gravedad inferior o igual a **1**.

2.12 Gráfica de frecuencia acumulada

F_i es una cantidad importante porque nos permite definir la acumulación de probabilidades hasta niveles intermedios.

La probabilidad de un nivel intermedio x ($i \leq x < i+1$) es solo la acumulación hasta el nivel inferior $F_x = F_i$.

F_x es por lo tanto una función de rango **continuo**. Podemos dibujarla con respecto a los resultados.



Por lo tanto, podemos decir que el **67 %** de los pacientes tenían misofonía hasta gravedad 2.3, aunque 2.3 no es un resultado observado.

2.13 Variables numéricas

El resultado de un experimento aleatorio puede producir un número. Si el número es **discreto**, podemos generar una tabla de frecuencias, con frecuencias absolutas, relativas y acumulativas, e ilustrarlas con gráficos de barras, de sectores y acumulativos.

Cuando el número es **continuo** las frecuencias no son útiles, lo más probable es que observemos o no un número continuo en particular.

Ejemplo (misofonía)

Los investigadores se preguntaron si la convexidad de la mandíbula afectaría la

gravedad de la misofonía. La hipótesis científica es que el ángulo de convexidad de la mandíbula puede influir en el oído y su sensibilidad. Estos son los resultados de la convexidad de la mandíbula (grados) para cada paciente:

```
## [1] 7.97 18.23 12.27 7.81 9.81 13.50 19.30 7.70 12.30 7.90 12.60 19.00
## [13] 7.27 14.00 5.40 8.00 11.20 7.75 7.94 16.69 7.62 7.02 7.00 19.20
## [25] 7.96 14.70 7.24 7.80 7.90 4.70 4.40 14.00 14.40 16.00 1.40 9.76
## [37] 7.90 7.90 7.40 6.30 7.76 7.30 7.00 11.23 16.00 7.90 7.29 6.91
## [49] 7.10 13.40 11.60 -1.00 6.00 7.82 4.80 11.00 9.00 11.50 16.00 15.00
## [61] 1.40 16.80 7.70 16.14 7.12 -1.00 17.00 9.26 18.70 3.40 21.30 7.50
## [73] 6.03 7.50 19.00 19.01 8.10 7.80 6.10 15.26 7.95 18.00 4.60 15.00
## [85] 7.50 8.00 16.80 8.54 7.00 18.30 7.80 16.00 14.00 12.30 11.40 8.50
## [97] 7.00 7.96 17.60 10.00 3.50 6.70 17.00 20.26 6.64 1.80 7.02 2.46
## [109] 19.00 17.86 6.10 6.64 12.00 6.60 8.70 14.05 7.20 19.70 7.70 6.02
## [121] 2.50 19.00 6.80
```

2.14 Transformando datos continuos

Como los resultados continuos no se pueden contar (de manera informativa), los transformamos en variables categóricas ordenadas.

- 1) Primero cubrimos el rango de las observaciones en intervalos regulares del mismo tamaño (contenedores)

```
## [1] "[-1.02,3.46]" "(3.46,7.92]" "(7.92,12.4]" "(12.4,16.8]" "(16.8,21.3]"
```

- 2) Luego mapeamos cada observación a su intervalo: creando una variable categórica **ordenada**; en este caso con 5 resultados posibles

```
## [1] "(7.92,12.4]" "(16.8,21.3]" "(7.92,12.4]" "(3.46,7.92]" "(7.92,12.4]"
## [6] "(12.4,16.8]" "(16.8,21.3]" "(3.46,7.92]" "(7.92,12.4]" "(3.46,7.92]"
## [11] "(12.4,16.8]" "(16.8,21.3]" "(3.46,7.92]" "(12.4,16.8]" "(3.46,7.92]"
## [16] "(7.92,12.4]" "(7.92,12.4]" "(3.46,7.92]" "(7.92,12.4]" "(12.4,16.8]"
## [21] "(3.46,7.92]" "(3.46,7.92]" "(3.46,7.92]" "(16.8,21.3]" "(7.92,12.4]"
## [26] "(12.4,16.8]" "(3.46,7.92]" "(3.46,7.92]" "(3.46,7.92]" "(3.46,7.92]"
## [31] "(3.46,7.92]" "(12.4,16.8]" "(12.4,16.8]" "(12.4,16.8]" "[-1.02,3.46]"
## [36] "(7.92,12.4]" "(3.46,7.92]" "(3.46,7.92]" "(3.46,7.92]" "(3.46,7.92]"
## [41] "(3.46,7.92]" "(3.46,7.92]" "(3.46,7.92]" "(7.92,12.4]" "(12.4,16.8]"
## [46] "(3.46,7.92]" "(3.46,7.92]" "(3.46,7.92]" "(3.46,7.92]" "(12.4,16.8]"
## [51] "(7.92,12.4]" "[-1.02,3.46]" "(3.46,7.92]" "(3.46,7.92]" "(3.46,7.92]"
## [56] "(7.92,12.4]" "(7.92,12.4]" "(7.92,12.4]" "(12.4,16.8]" "(12.4,16.8]"
## [61] "[-1.02,3.46]" "(12.4,16.8]" "(3.46,7.92]" "(12.4,16.8]" "(3.46,7.92]"
## [66] "[-1.02,3.46]" "(16.8,21.3]" "(7.92,12.4]" "(16.8,21.3]" "[-1.02,3.46]"
## [71] "(16.8,21.3]" "(3.46,7.92]" "(3.46,7.92]" "(3.46,7.92]" "(16.8,21.3]"
## [76] "(16.8,21.3]" "(7.92,12.4]" "(3.46,7.92]" "(3.46,7.92]" "(12.4,16.8]"
## [81] "(7.92,12.4]" "(16.8,21.3]" "(3.46,7.92]" "(12.4,16.8]" "(3.46,7.92]"
## [86] "(7.92,12.4]" "(12.4,16.8]" "(7.92,12.4]" "(3.46,7.92]" "(16.8,21.3]"
```


2.15. TABLA DE FRECUENCIAS PARA UNA VARIABLE CONTINUA 17

```
## [91] "(3.46,7.92]" "(12.4,16.8]" "(12.4,16.8]" "(7.92,12.4]" "(7.92,12.4]"
## [96] "(7.92,12.4]" "(3.46,7.92]" "(7.92,12.4]" "(16.8,21.3]" "(7.92,12.4]"
## [101] "(3.46,7.92]" "(3.46,7.92]" "(16.8,21.3]" "(16.8,21.3]" "(3.46,7.92]"
## [106] "[-1.02,3.46]" "(3.46,7.92]" "[-1.02,3.46]" "(16.8,21.3]" "(16.8,21.3]"
## [111] "(3.46,7.92]" "(3.46,7.92]" "(7.92,12.4]" "(3.46,7.92]" "(7.92,12.4]"
## [116] "(12.4,16.8]" "(3.46,7.92]" "(16.8,21.3]" "(3.46,7.92]" "(3.46,7.92]"
## [121] "[-1.02,3.46]" "(16.8,21.3]" "(3.46,7.92]"
```

Por tanto, en lugar de decir que el primer paciente tenía un ángulo de convexidad de 7.97, decimos que su ángulo estaba entre el intervalo (o **bin**) (7.92, 12.4].

Ningún otro paciente tenía un ángulo de 7.97, pero muchos tenían ángulos entre (7.92, 12.4].

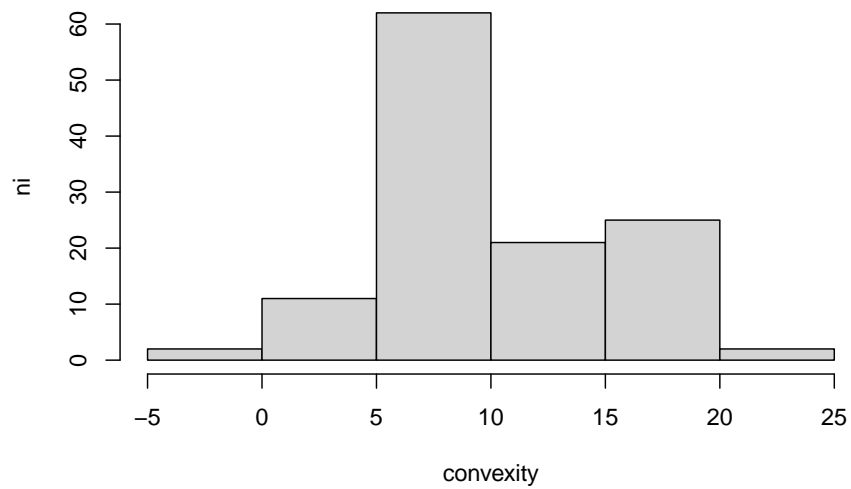
2.15 Tabla de frecuencias para una variable continua

Para una partición regular dada del intervalo de resultados en intervalos, podemos producir una tabla de frecuencias como antes

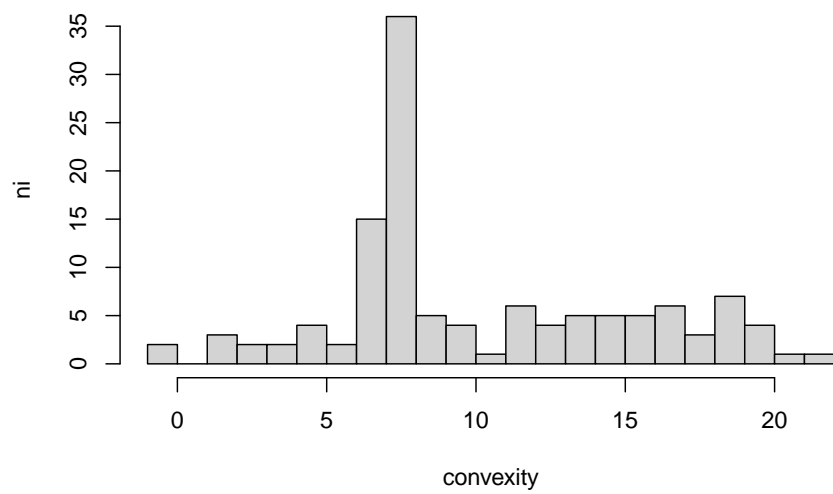
##	outcome	ni	fi	Ni	Fi
## 1	[-1.02,3.46]	8	0.06504065	8	0.06504065
## 2	(3.46,7.92]	51	0.41463415	59	0.47967480
## 3	(7.92,12.4]	26	0.21138211	85	0.69105691
## 4	(12.4,16.8]	20	0.16260163	105	0.85365854
## 5	(16.8,21.3]	18	0.14634146	123	1.00000000

2.16 Histograma

El histograma es la gráfica de n_i o f_i Vs los resultados en intervalos (bins). El histograma depende del tamaño de los bins.



Este es un histograma con 20 bins.



Vemos que la mayoría de las personas tienen ángulos dentro de $(7, 8]$

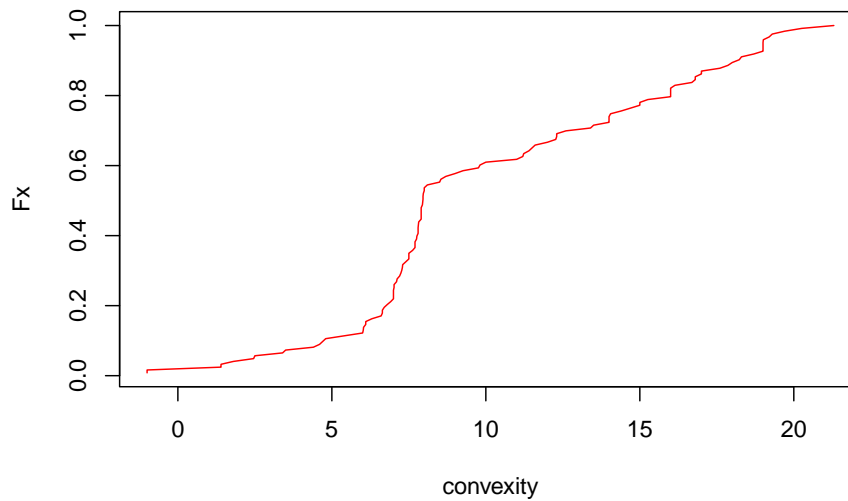
2.17 Gráfica de frecuencia acumulada

También podemos graficar F_x contra los resultados. Como F_x es de rango continuo, podemos ordenar las observaciones ($x_1 < \dots x_j < x_{j+1} < x_n$) y por lo tanto

$$F_x = \frac{k}{n}$$

para $x_k \leq x < x_{k+1}$.

F_x se conoce como la **distribución** de los datos. F_x no depende del tamaño del bin. Sin embargo, su **resolución** depende de la cantidad de datos.



2.18 Estadísticas de resumen

Las estadísticas de resumen son números calculados a partir de los datos que nos dicen características importantes de las variables numéricas (discretas o continuas).

Por ejemplo, tenemos estadísticas que describen los valores extremos:

- **mínimo:** el resultado mínimo observado
- **máximo:** el resultado máximo observado

2.19 Promedio (media muestral)

Una estadística importante que describe el valor central de los resultados (dónde esperar la mayoría de las observaciones) es el **promedio**

$$\bar{x} = \frac{1}{N} \sum_{j=1..N} x_j$$

donde x_j es la **observación** j de un total de N .

Ejemplo (Misofonía)

La convexidad promedio se puede calcular directamente a partir de las **observaciones**

$$\begin{aligned} \bar{x} &= \frac{1}{N} \sum_j x_j \\ &= \frac{1}{N} (7.97 + 18.23 + 12.27 \dots + 6.80) = 10.19894 \end{aligned}$$

Para variables **categoricamente ordenadas**, podemos usar las frecuencias relativas para calcular el promedio

$$\begin{aligned} \bar{x} &= \frac{1}{N} \sum_{i=1..N} x_j = \frac{1}{N} \sum_{i=1..M} x_i * n_i \\ &= \sum_{i=1..M} x_i * f_i \end{aligned}$$

donde pasamos de sumar N **observaciones** a sumar M **resultados**.

La forma $\bar{x} = \sum_{i=1..M} x_i f_i$ muestra que el promedio es el **centro de gravedad** de los resultados. Como si cada resultado tuviera una densidad de masa dada por f_i .

Ejemplo (Misofonía)

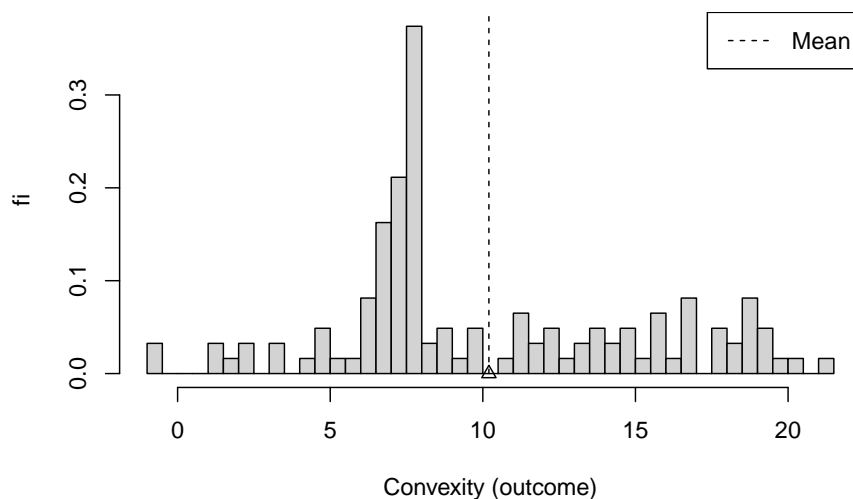
La **severidad** promedio de la misofonía en el estudio se puede calcular a partir de las frecuencias relativas de los **resultados**

```
## outcome ni      fi
## 1      0 41 0.3333333
## 2      1  5 0.04065041
## 3      2 37 0.30081301
## 4      3 31 0.25203252
## 5      4  9 0.07317073
```

$$\bar{x} = 0 * f_0 + 1 * f_1 + 2 * f_2 + 3 * f_3 + 4 * f_4 = 1.691057$$

2.20 Promedio

El promedio es también el centro de gravedad de las variables continuas. Ese es el punto donde las frecuencias reativas se equilibran.



2.21 mediana

Otra medida de centralidad es la mediana. La mediana x_m , o $q_{0.5}$, es el valor por debajo del cual encontramos la mitad de las observaciones. Cuando ordenamos las observaciones $x_1 < \dots < x_j < x_{j+1} < x_N$, las contamos hasta encontrar la mitad de ellas. x_m es tal que

$$\sum_{i \leq m} 1 = \frac{N}{2}$$

Ejemplo (Misofonía)

Si ordenamos los ángulos de convexidad, vemos que 62 observaciones (individuos) ($N/2 \sim 123/2$) están por debajo de 7.96. La **convexidad mediana** es por lo tanto $q_{0.5} = x_{62} = 7.96$

##	[1]	-1.00	-1.00	1.40	1.40	1.80	2.46	2.50	3.40	3.50	4.40	4.60	4.70
##	[13]	4.80	5.40	6.00	6.02	6.03	6.10	6.10	6.30	6.60	6.64	6.64	6.70
##	[25]	6.80	6.91	7.00	7.00	7.00	7.00	7.02	7.02	7.10	7.12	7.20	7.24
##	[37]	7.27	7.29	7.30	7.40	7.50	7.50	7.50	7.62	7.70	7.70	7.70	7.75
##	[49]	7.76	7.80	7.80	7.80	7.81	7.82	7.90	7.90	7.90	7.90	7.90	7.94

```
## [61] 7.95 7.96

## [1] 7.96 7.97 8.00 8.00 8.10 8.50 8.54 8.70 9.00 9.26 9.76 9.81
## [13] 10.00 11.00 11.20 11.23 11.40 11.50 11.60 12.00 12.27 12.30 12.30 12.60
## [25] 13.40 13.50 14.00 14.00 14.00 14.05 14.40 14.70 15.00 15.00 15.26 16.00
## [37] 16.00 16.00 16.00 16.14 16.69 16.80 16.80 17.00 17.00 17.60 17.86 18.00
## [49] 18.23 18.30 18.70 19.00 19.00 19.00 19.00 19.01 19.20 19.30 19.70 20.26
## [61] 21.30

## [1] 7.96
```

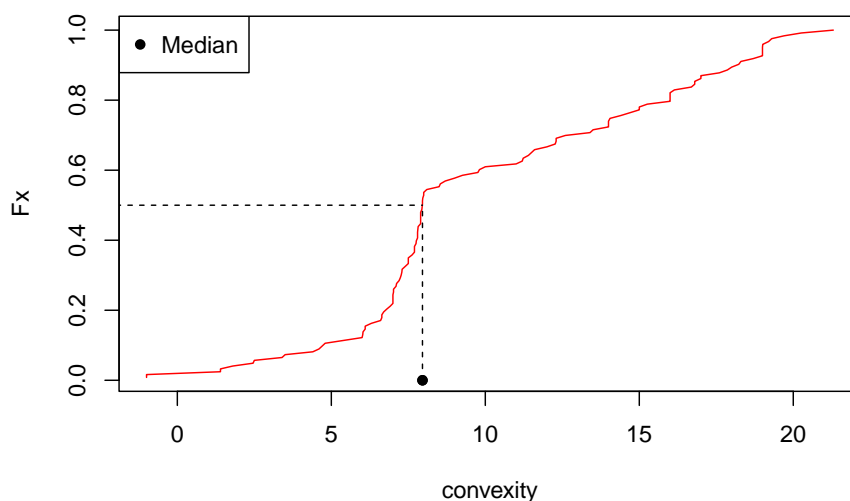
En términos de frecuencias, $q_{0.5}$ hace que la frecuencia acumulada F_x sea igual a 0.5

$$\sum_{i=0, \dots, m} f_i = F_{q_{0.5}} = 0.5$$

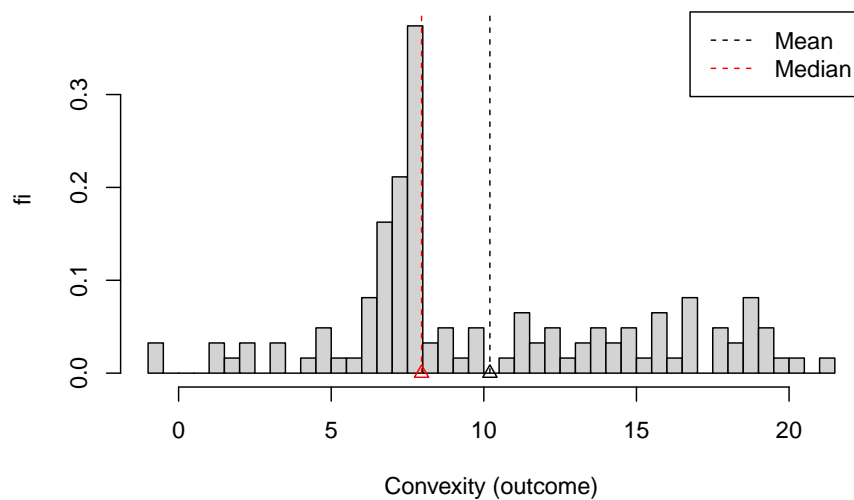
o

$$q_{0.5} = F^{-1}(0.5)$$

En el gráfico de distribución, la mediana es el valor de x en el que se encuentra la mitad del máximo de F .



El promedio y la mediana no siempre son iguales.

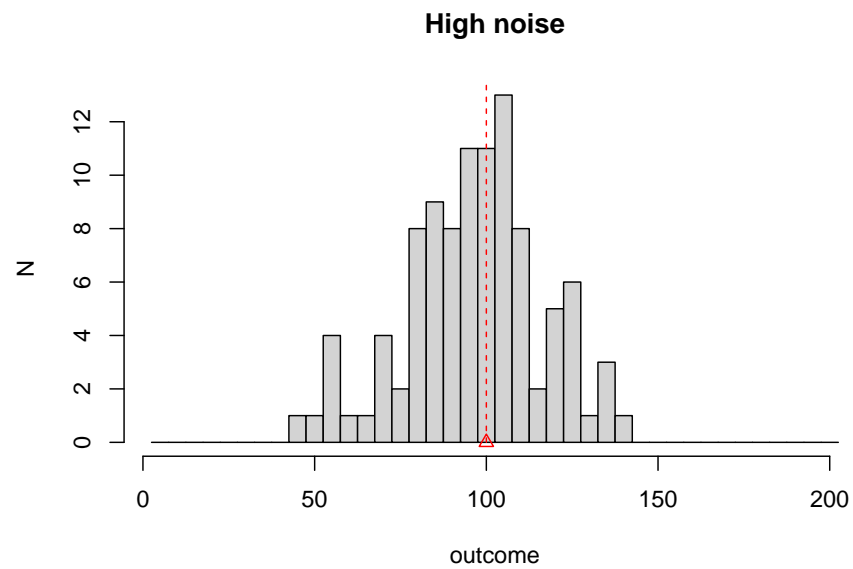
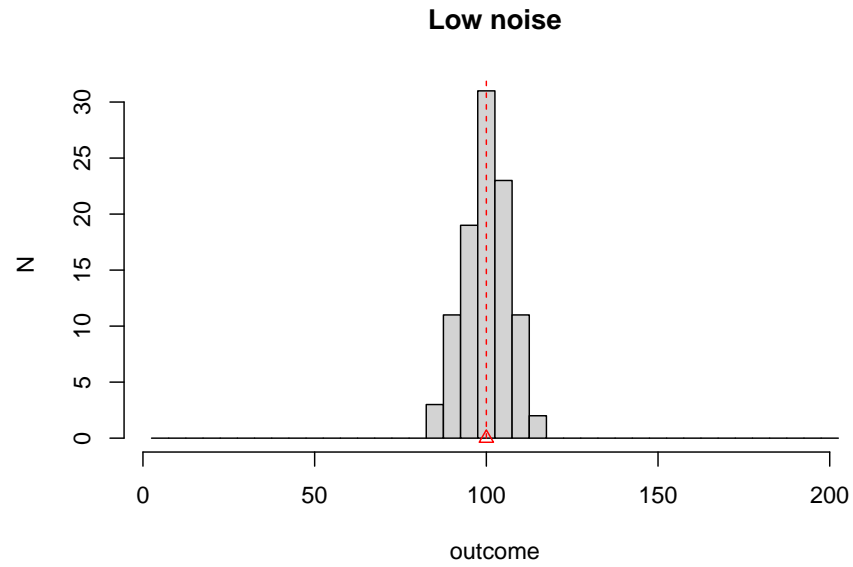


2.22 Dispersión

Otras estadísticas de resumen importantes de las observaciones son las de **dispersión**.

Muchos experimentos pueden compartir su media, pero difieren en cuán **dispersos** son los valores.

La dispersión de las observaciones es una medida del **ruido**.



2.23 Variación de la muestra

La dispersión sobre la media se mide con la varianza muestral

$$s^2 = \frac{1}{N-1} \sum_{j=1..N} (x_j - \bar{x})^2$$

Este número, mide la distancia cuadrada promedio de las **observaciones** al promedio. La razón de $N-1$ se explicará cuando hablemos de inferencia, cuando estudiemos la dispersión de \bar{x} , además de la dispersión de las observaciones.

En términos de las frecuencias de las variables **categorías y ordenadas**

$$s^2 = \frac{N}{N-1} \sum_{i=1..M} (x_i - \bar{x})^2 f_i$$

s^2 se puede considerar como el **momento de inercia** de las observaciones.

La raíz cuadrada de la varianza de la muestra se denomina **desviación estándar** s .

Ejemplo (Misofonía)

La desviación estándar del ángulo de convexidad es

$$s = [\frac{1}{123-1} ((7.97 - 10.19894)^2 + (18.23 - 10.19894)^2 + (12.27 - 10.19894)^2 + \dots)]^{1/2} = 5.086707$$

La convexidad de la mandíbula se desvía de su media en 5.086707.

2.24 Rango intercuartílico (IQR)

La dispersión de los datos también se puede medir con respecto a la mediana usando el **rango intercuartílico**:

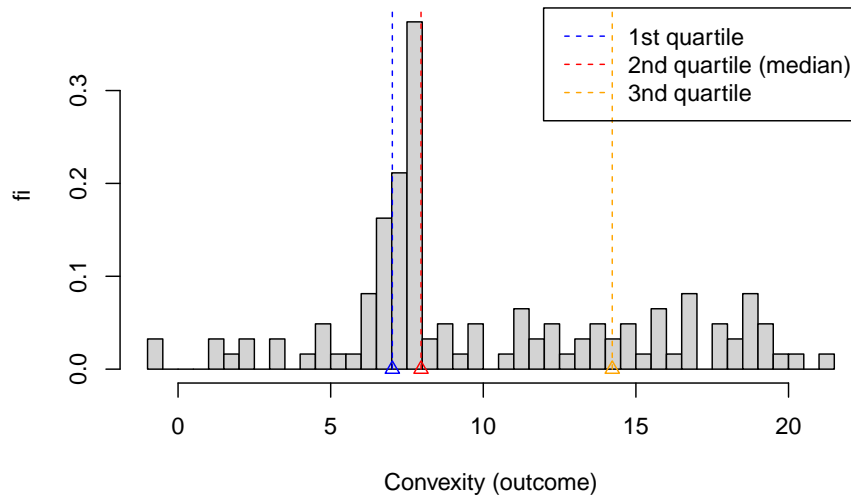
- 1) Definimos el **primer** cuartil como el valor x_m que hace que la frecuencia acumulada $F_{q_{0.25}}$ sea igual a 0.25 (x donde hemos acumulado una cuarta parte de las observaciones)

$$F_{q_{0.25}} = 0.25$$

- 1) Definimos el **tercer** cuartil como el valor x_m que hace que la frecuencia acumulada $F_{q_{0.75}}$ sea igual a 0.75 (x donde hemos acumulado tres cuartos de observaciones)

$$F_{q_{0.75}} = 0.75$$

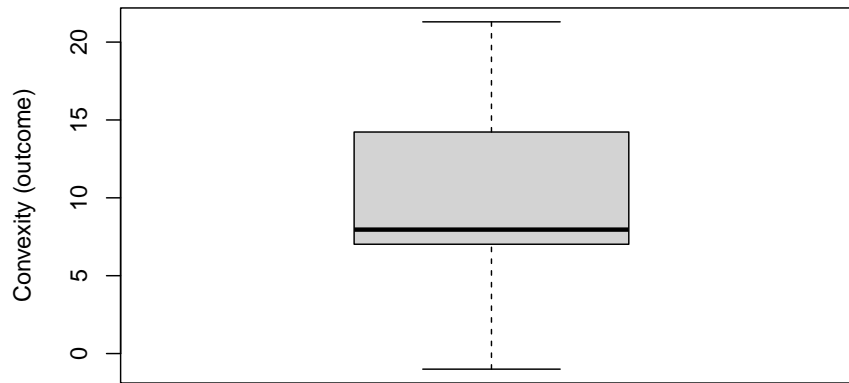
- 3) El **rango intercuartílico** (IQR) es $IQR = q_{0.75} - q_{0.25}$. Esa es la distancia entre el tercer y el primer cuartil y captura el 50% central de las observaciones



2.25 Diagrama de caja

El rango intercuartílico, la mediana y los 5% y 95% de los datos se pueden visualizar en un **diagrama de caja**.

En el diagrama de caja, los valores de los resultados están en el eje y. El IQR es la caja, la mediana es la línea del medio y los bigotes marcan los 5% y 95% de los datos.



2.26 Preguntas

1) En el siguiente diagrama de caja, el primer cuartil y el segundo cuartil de los datos son:

a: $(-1.00, 21.30)$; **b:** $(-1.00, 7.02)$; **c:** $(7.02, 7.96)$; **d:** $(7.02, 14.22)$

2) La principal desventaja de un histograma es que:

a: Depende del tamaño del bin; **b:** No se puede utilizar para variables categóricas; **c:** No se puede usar cuando el tamaño del bin es pequeño; **d:** Se usa solo para frecuencias relativas;

3) Si las frecuencias acumuladas relativas de un experimento aleatorio con resultados $\{1, 2, 3, 4\}$ son: $F(1) = 0.15$, $F(2) = 0.60$, $F(3) = 0.85$, $F(4) = 1$.

Entonces la frecuencia relativa para el resultado 3 es

a: 0.15; **b:** 0.85; **c:** 0.45; **d:** 0.25

4) En una muestra de tamaño 10 de un experimento aleatorio obtuvimos los siguientes datos:

8, 3, 3, 7, 3, 6, 5, 10, 3, 8.

El primer cuartil de los datos es:

a: 3.5; **b:** 4; **c:** 5; **d:** 3

5) Imaginemos que recopilamos datos para dos cantidades que no son mutuamente excluyentes, por ejemplo, el sexo y la nacionalidad de los pasajeros de un vuelo. Si queremos hacer un solo gráfico circular para los datos, ¿cuál de estas afirmaciones es verdadera?

a: Solo podemos hacer un gráfico circular de nacionalidad porque tiene más de dos resultados posibles; **b:** Podemos hacer un gráfico circular para una variable nueva que marca el sexo **y** la nacionalidad; **c:** Podemos hacer un gráfico circular para la variable sexo o la variable nacionalidad; **d:** Solo podemos elegir si hacemos un gráfico circular para el sexo o un gráfico circular para la nacionalidad.

2.27 Ejercicios

2.27.0.1 Ejercicio 1

Hemos realizado un experimento 8 veces con los siguientes resultados

```
## [1] 3 3 10 2 6 11 5 4
```

Responde las siguientes cuestiones:

- Calcula las frecuencias relativas de cada resultado.
- Calcula las frecuencias acumuladas de cada resultado.
- ¿Cuál es el promedio de las observaciones?
- ¿Cuál es la mediana?
- ¿Cuál es el tercer cuartil?
- ¿Cuál es el primer cuartil?

2.27.0.2 Ejercicio 2

Hemos realizado un experimento 10 veces con los siguientes resultados

```
## [1] 2.875775 7.883051 4.089769 8.830174 9.404673 0.455565 5.281055 8.924190
## [9] 5.514350 4.566147
```

Considera 10 bins de tamaño 1: $[0,1]$, $(1,2]$... $(9,10]$.

Responde las siguientes cuestiones:

- Calcula las frecuencias relativas de cada resultado y dibuja el histograma
- Calcula las frecuencias acumulativas de cada resultado y dibuja la gráfica acumulativa.
- Dibuja un diagrama de caja .

Chapter 3

Probabilidad

En este capítulo introduciremos el concepto de probabilidad a partir de frecuencias relativas.

Definiremos los eventos como los elementos sobre los que se aplica la probabilidad. Los eventos compuestos se definirán usando álgebra de conjuntos.

Luego discutiremos el concepto de probabilidad condicional derivado de la probabilidad conjunta de dos eventos.

3.1 Experimentos aleatorios

Recordemos el objetivo básico de la estadística. La estadística se ocupa de los datos que se presentan en forma de observaciones.

- Una **observación** es la adquisición de un número o una característica de un experimento

Las observaciones son realizaciones de **resultados**.

- Un **resultado** es una posible observación que es el resultado de un experimento.

Al realizar experimentos, a menudo obtenemos resultados diferentes. La descripción de la variabilidad de los resultados es uno de los objetivos de la estadística.

- Un **experimento aleatorio** es un experimento que da diferentes resultados cuando se repite de la misma manera.

La pregunta filosófica detrás es ¿Cómo podemos conocer algo si cada vez que lo miramos cambia?

3.2 Probabilidad de medición

Nos gustaría tener una medida para el resultado de un experimento aleatorio que nos diga **cuán seguros** estamos de observar el resultado cuando realicemos un **futuro** experimento aleatorio.

Llamaremos a esta medida la probabilidad del resultado y le asignaremos valores:

- 0, cuando estamos seguros de que la observación **no** ocurrirá.
- 1, cuando estamos seguros de que la observación sucederá.

3.3 Probabilidad clásica

Siempre que un experimento aleatorio tenga M resultados posibles que son todos **igualmente probables**, la probabilidad de cada resultado i es

$$P_i = \frac{1}{M}$$

.

La probabilidad clásica fue defendida por Laplace (1814).

Dado que cada resultado es **igualmente probable** en este tipo de experimento, declaramos una completa ignorancia y lo mejor que podemos hacer es distribuir equitativamente la misma probabilidad para cada resultado.

- No observamos P_i
- Deducimos P_i de nuestra razón y no necesitamos realizar ningún experimento para conocerla.

Ejemplo (dado):

¿Cuál es la probabilidad de que obtengamos 2 en el lanzamiento de un dado?

$$P_2 = 1/6 = 0.166666.$$

3.4 Frecuencias relativas

¿Qué sucede con los experimentos aleatorios cuyos posibles resultados **no** son igualmente probables?

¿Cómo podemos entonces definir las probabilidades de los resultados?

Ejemplo (experimento aleatorio)

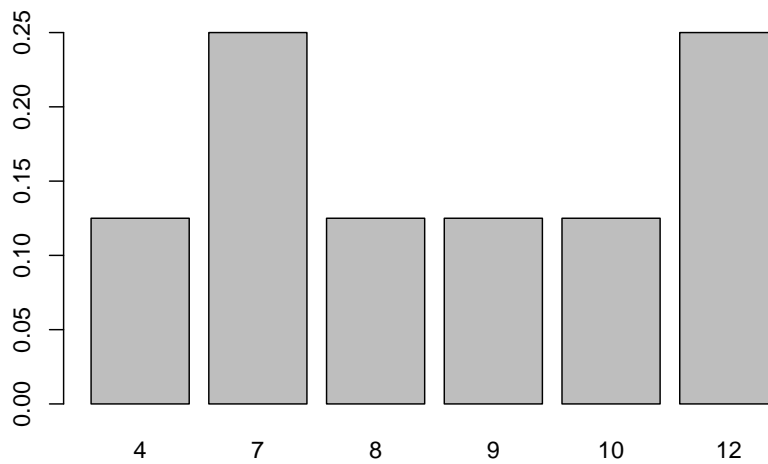
Imaginemos que repetimos un experimento aleatorio 8 veces y obtenemos las siguientes observaciones

8 4 12 7 10 7 9 12

- ¿Qué tan seguro estamos de obtener el resultado 12 en la siguiente observación?

La tabla de frecuencias es

##	outcome	n_i	f_i
## 1	4	1	0.125
## 2	7	2	0.250
## 3	8	1	0.125
## 4	9	1	0.125
## 5	10	1	0.125
## 6	12	2	0.250



La **frecuencia relativa** $f_i = \frac{n_i}{N}$ parece una medida de probabilidad razonable porque

- es un número entre 0 y 1.
- mide la proporción del total de observaciones que observamos de un resultado particular.

Como $f_{12} = 0.25$ entonces estaríamos un cuarto seguros, una de cada 4 observaciones, de obtener 12.

Pregunta: ¿Qué tan bueno es f_i como medida de certeza del resultado i ?

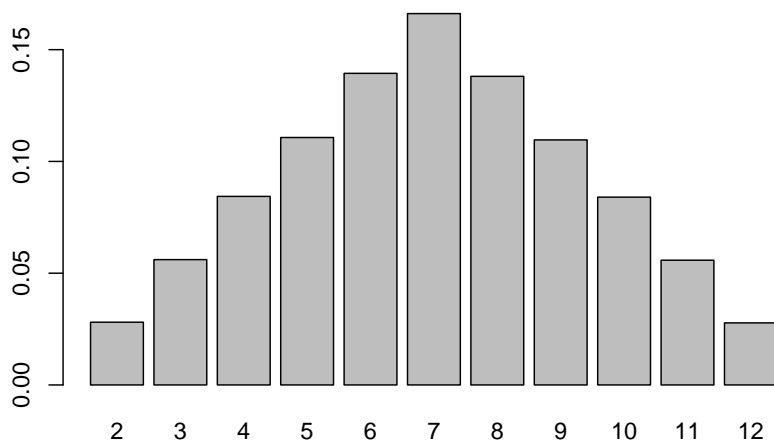
Ejemplo (experimento aleatorio con mas repeticiones)

Digamos que repetimos el experimento 100000 veces más:

La tabla de frecuencias es ahora

##	outcome	ni	fi
## 1	2	2807	0.02807
## 2	3	5607	0.05607
## 3	4	8435	0.08435
## 4	5	11070	0.11070
## 5	6	13940	0.13940
## 6	7	16613	0.16613
## 7	8	13806	0.13806
## 8	9	10962	0.10962
## 9	10	8402	0.08402
## 10	11	5581	0.05581
## 11	12	2777	0.02777

y el gráfico de barras es



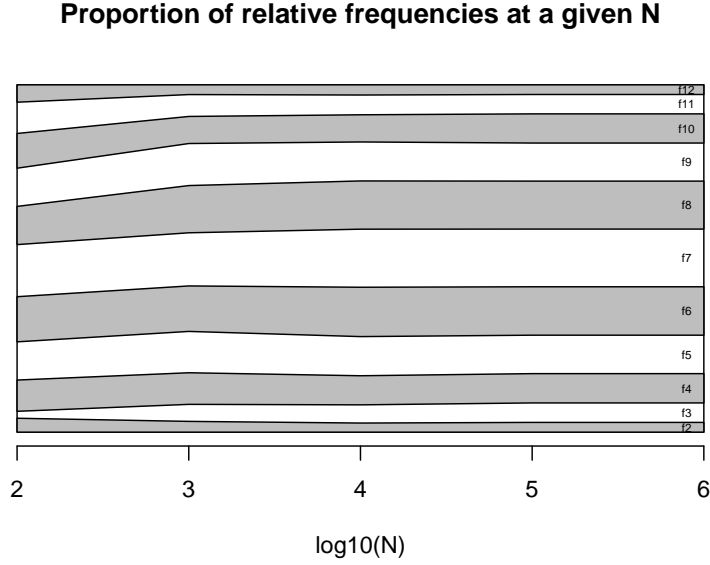
Aparecieron nuevos resultados y f_{12} ahora es solo 0.027, y entonces estamos sólo un $\sim 3\%$ seguros de obtener 12 en el próximo experimento. Las probabilidades medidas por f_i cambian con N .

3.5 Frecuencias relativas en el infinito

Una observación crucial es que si medimos las probabilidades de f_i en valores crecientes de N ¡convergen!

En este gráfico cada sección vertical da la frecuencia relativa de cada observación. Vemos que después de $N = 1000$ ($\log_{10}(N) = 3$) las proporciones apenas varían

con mas N .



Encontramos que cada una de las frecuencias relativas f_i converge a un valor constante

$$\lim_{N \rightarrow \infty} f_i = P_i$$

3.6 Probabilidad frecuentista

Llamamos **Probabilidad** P_i al límite cuando $N \rightarrow \infty$ de la **frecuencia relativa** de observar el resultado i en un experimento aleatorio.

Defendida por Venn (1876), la definición frecuentista de probabilidad se deriva de datos/experiencia (empírica).

- No observamos P_i , observamos f_i
- **Estimamos** P_i con f_i (normalmente cuando N es grande), escribimos:

$$\hat{P}_i = f_i$$

Similar a la relación entre **observación** y **resultado**, tenemos la relación entre **frecuencia relativa** y **probabilidad** como un valor concreto de una cantidad abstracta.

3.7 Probabilidades clásicas y frecuentistas

Tenemos situaciones en las que se puede usar la probabilidad clásica para encontrar el límite de frecuencias relativas.

- Si los resultados son **igualmente probables**, la probabilidad clásica nos da el límite:

$$P_i = \lim_{N \rightarrow \infty} \frac{n_i}{N} = \frac{1}{M}$$

- Si los resultados en los que estamos interesados pueden derivarse de otros resultados **igualmente probables**; Veremos más sobre esto cuando estudiemos los modelos de probabilidad.

Ejemplo (suma de dos dados)

Nuestro ejemplo anterior se basa en la **suma de dos dados**. Si bien realizamos el experimento muchas veces, anotamos los resultados y calculamos las **frecuencias relativas**, podemos conocer el valor exacto de probabilidad.

Esta probabilidad **se deduce** del hecho de que el resultado de cada dado es **igualmente probable**. A partir de esta suposición, podemos encontrar que (Ejercicio 1)

$$P_i = \begin{cases} \frac{i-1}{36}, & i \in \{2, 3, 4, 5, 6, 7\} \\ \frac{13-i}{36}, & i \in \{8, 9, 10, 11, 12\} \end{cases}$$

La motivación de la definición frecuentista es **empírica** (datos) mientras que la de la definición clásica es **racional** (modelos). A menudo combinamos ambos enfoques (inferencia y deducción) para conocer las probabilidades de nuestro experimento aleatorio.

3.8 Definición de probabilidad

Una probabilidad es un número que se asigna a cada resultado posible de un experimento aleatorio y satisface las siguientes propiedades o **axiomas**:

- 1) cuando los resultados E_1 y E_2 son mutuamente excluyentes; es decir, solo uno de ellos puede ocurrir, entonces la probabilidad de observar E_1 o E_2 , escrito como $E_1 \cup E_2$, es su suma:

$$P(E_1 \cup E_2) = P(E_1) + P(E_2)$$

- 2) cuando S es el conjunto de todos los resultados posibles, entonces su probabilidad es 1 (al menos se observa algo):

$$P(S) = 1$$

- 3) La probabilidad de cualquier resultado está entre 0 y 1

$$P(E) \in [0, 1]$$

Propuesto por Kolmogorov's hace menos de 100 años (1933)

3.9 Tabla de probabilidades

Las propiedades de Kolmogorov son las reglas básicas para construir una **tabla de probabilidad**, de manera similar a la tabla de frecuencia relativa.

Ejemplo (Dado)

La tabla de probabilidad para el lanzamiento de un dado

resultado	Probabilidad
1	1/6
2	1/6
3	1/6
4	1/6
5	1/6
6	1/6
$P(1 \cup 2 \cup \dots \cup 6)$	1

Verifiquemos los axiomas:

- 1) Donde $1 \cup 2$ es, por ejemplo, el **evento** de lanzar un 1 o un 2. Entonces

$$P(1 \cup 2) = P(1) + P(2) = 2/6$$

- 2) Como $S = \{1, 2, 3, 4, 5, 6\}$ se compone de resultados **mutuamente excluyentes**, entonces

$$P(S) = P(1 \cup 2 \cup \dots \cup 6) = P(1) + P(2) + \dots + P(n) = 1$$

- 3) Las probabilidades de cada uno de resultados están entre 0 y 1.

3.10 Espacio muestral

El conjunto de todos los resultados posibles de un experimento aleatorio se denomina **espacio muestral** y se denota como S .

El espacio muestral puede estar formado por resultados categóricos o numéricos.

Por ejemplo:

- temperatura humana: $S = (36, 42)$ grados Celsius.

- niveles de azúcar en humanos: $S = (70 - 80)mg/dL$
- el tamaño de un tornillo de una línea de producción: $S = (70 - 72)mm$
- número de correos electrónicos recibidos en una hora: $S = \{1, \dots, \infty\}$
- el lanzamiento de un dado: $S = \{1, 2, 3, 4, 5, 6\}$

3.11 Eventos

Un **evento** A es un **subconjunto** del espacio muestral. Es una **colección** de resultados.

Ejemplos de eventos:

- El evento de una temperatura saludable: $A = 37 - 38$ grados Celsius
- El evento de producir un tornillo con un tamaño: $A = 71.5mm$
- El evento de recibir más de 4 emails en una hora: $A = \{4, \infty\}$
- El evento de obtener un número menor o igual a 3 en la tirada de a dice:
 $A = \{1, 2, 3\}$

Un evento se refiere a un posible conjunto de **resultados**.

3.12 Álgebra de eventos

Para dos eventos A y B , podemos construir los siguientes eventos derivados utilizando las operaciones básicas de conjuntos:

- Complemento A' : el evento de **no** A
- Unión $A \cup B$: el evento de A **o** B
- Intersección $A \cap B$: el evento de A **y** B

Ejemplo (dado)

Lancemos un dado y veamos los eventos (conjunto de resultados):

- un número menor o igual a tres $A : \{1, 2, 3\}$
- un número par $B : \{2, 4, 6\}$

Veamos como podemos construir nuevos eventos con las operaciones de conjuntos:

- un número no menor de tres: $A' : \{4, 5, 6\}$
- un número menor o igual a tres **o** par: $A \cup B : \{1, 2, 3, 4, 6\}$
- un número menor o igual a tres **y** par $A \cap B : \{2\}$

3.13 Resultados mutuamente excluyentes

Los resultados como tirar 1 y 2 en un dado son eventos que no pueden ocurrir al mismo tiempo. Decimos que son **mutuamente excluyentes**.

En general, dos eventos denotados como E_1 y E_2 son mutuamente excluyentes cuando

$$E_1 \cap E_2 = \emptyset$$

Ejemplos:

- El resultado de tener una gravedad de misofonía de 1 y una gravedad de 4.
- Los resultados de obtener 12 y 5 al sumar el lanzamiento de dos dados.

De acuerdo con las propiedades de Kolmogorov, solo los resultados **mutuamente excluyentes** se pueden organizar en **tablas de probabilidad**, como en las tablas de frecuencias relativas.

3.14 Probabilidades conjuntas

La **probabilidad conjunta** de A y B es la probabilidad de A y B . Eso es

$$P(A \cap B)$$

o $P(A, B)$.

Para escribir probabilidades conjuntas de eventos no mutuamente excluyentes ($A \cap B \neq \emptyset$) en una tabla de probabilidad, notamos que siempre podemos descomponer el espacio muestral en conjuntos **mutuamente excluyentes** que involucran las intersecciones:

$$S = \{A \cap B, A \cap B', A' \cap B, A' \cap B'\}$$

Consideremos el diagrama de Ven para el ejemplo donde A es el evento que corresponde a sacar número menor o igual que 3 y B corresponde a un número par:

Las **marginales** de A y B son la probabilidad de A y la probabilidad de B , respectivamente:

- $P(A) = P(A \cap B') + P(A \cap B) = 2/6 + 1/6 = 3/6$
- $P(B) = P(A' \cap B) + P(A \cap B) = 2/6 + 1/6 = 3/6$

Podemos ahora escribir la **tabla de probabilidad** para las probabilidades conjuntas

Resultado	Probabilidad
$(A \cap B)$	$P(A \cap B) = 1/6$
$(A \cap B')$	$P(A \cap B') = 2/6$
$(A' \cap B)$	$P(A' \cap B) = 2/6$
$(A' \cap B')$	$P(A' \cap B') = 1/6$

Resultado	Probabilidad
suma	1

Cada resultado tiene *dos* valores (uno para la característica del tipo A y otro para el tipo B)

3.15 Tabla de contingencia

La tabla de probabilidad conjunta también se puede escribir en una **tabla de contingencia**

	B	B'	suma
A	$P(A \cap B)$	$P(A \cap B')$	$P(A)$
A'	$P(A' \cap B)$	$P(A' \cap B')$	$P(A')$
suma	$P(B)$	$P(B')$	1

Donde las marginales son las sumas en las márgenes de la tabla, por ejemplo:

- $P(A) = P(A \cap B') + P(A \cap B)$
- $P(B) = P(A' \cap B) + P(A \cap B)$

En nuestro ejemplo, la tabla de contingencia es

	B	B'	suma
A	1/6	2/6	3/6
A'	2/6	1/6	3/6
suma	3/6	3/6	1

3.16 La regla de la suma:

La regla de la suma nos permite calcular la probabilidad de A o B , $P(A \cup B)$, en términos de la probabilidad de A y B , $P(A \cap B)$. Podemos hacer esto de tres maneras equivalentes:

- 1) Usando las marginales y la probabilidad conjunta

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

- 2) Usando solo probabilidades conjuntas

$$P(A \cup B) = P(A \cap B) + P(A \cap B') + P(A' \cap B)$$

3) Usando el complemento de la probabilidad conjunta

$$P(A \cup B) = 1 - P(A' \cap B')$$

Ejemplo (dado)

Tomemos los eventos $A : \{1, 2, 3\}$, sacar un número menor o igual que 3, y $B : \{2, 4, 6\}$, sacar un número par en el lanzamiento de un dado.

Por lo tanto:

$$1) P(A \cup B) = P(A) + P(B) - P(A \cap B) = 3/6 + 3/6 - 1/6 = 5/6$$

$$2) P(A \cup B) = P(A \cap B) + P(A \cap B') + P(A' \cap B) = 1/6 + 2/6 + 2/6 = 5/6$$

$$3) P(A \cup B) = 1 - P(A' \cap B') = 1 - 1/6 = 5/6$$

En la tabla de contingencia $P(A \cup B)$ corresponde a las casillas en negrita (o sea todas menos $1/6$ de abajo a la derecha)

	B	B'
A	1/6	2/6
A'	2/6	1/6

3.17 Preguntas

Recopilamos la edad y categoría de 100 deportistas en una competición

	<i>edad : junior</i>	<i>edad : senior</i>
<i>categoría : 1ra</i>	14	12
<i>categoría : 2a</i>	21	18
<i>categoría : 3a</i>	22	13

1) ¿Cuál es la probabilidad estimada de que un deportista sea de 2ª categoría y senior?

a: 18/100; **b:** 18/43; **c:** 18; **d:** 18/39

2) ¿Cuál es la probabilidad estimada de que el atleta no esté en la tercera categoría y sea senior?

a: 35/100; **b:** 30/100; **c:** 22/100; **d:** 13/100

3) ¿Cuál es la probabilidad marginal de la tercera categoría?

a: 13/100; **b:** 35/100; **c:** 22/100; **d:** 13/22

4) ¿Cuál es la probabilidad marginal de ser senior?

a: 13/100; b: 43/100; c: 43/57; d: 57/100

5) ¿Cuál es la probabilidad de ser senior o de tercera categoría?

a: 65/100; b: 86/100; c: 78/100; d: 13/100

3.18 Ejercicios

3.18.0.1 Probabilidad clásica: Ejercicio 1

- Escribe la tabla de **probabilidad conjunta** para los **resultados** de lanzar dos dados; en las filas escribe los resultados del primer dado y en las columnas los resultados del segundo dado.
- ¿Cuál es la probabilidad de sacar (3, 4)? (R:1/36)
- ¿Cuál es la probabilidad de tirar 3 y 4 con cualquiera de los dos dados? (R:2/36)
- ¿Cuál es la probabilidad de tirar 3 en el primer dado o 4 en el segundo? (A:11/36)
- ¿Cuál es la probabilidad de tirar 3 o 4 con cualquier dado? (R:20/36)
- Escribe la **tabla de probabilidad** para el resultado de la **suma** de dos dados. Supon que el resultado de cada dado es **igualmente probable**. Verifica que es:

$$P_i = \begin{cases} \frac{i-1}{36}, & i \in \{2, 3, 4, 5, 6, 7\} \\ \frac{13-i}{36}, & i \in \{8, 9, 10, 11, 12\} \end{cases}$$

3.18.0.2 Probabilidad frecuentista: Ejercicio 2

El resultado de un experimento aleatorio es medir la gravedad de la misofonía y el estado de depresión de un paciente.

- Gravedad de la misofonía: $x \in \{0, 1, 2, 3, 4\}$
- Depresión: $y \in \{0, 1\}$ (no:0, si:1)

Estos son los primeros 6 pacientes:

```
## Misofonia.dic depresion.dic
## 1             4             1
## 2             2             0
## 3             0             0
## 4             3             0
## 5             0             0
## 6             0             0
```

El estudio sobre un total de 123 pacientes mostró las frecuencias $n_{x,y}$ dadas en la tabla de contingencia:


```
##
##           Depression:0 Depression:1
## Misophonia:4           0           9
## Misophonia:3          25           6
## Misophonia:2          34           3
## Misophonia:1           5           0
## Misophonia:0          36           5
```

Supongamos que N es grande y que las frecuencias **estiman** las probabilidades $f_{x,y} = \hat{P}(X,Y)$, de tal forma que la siguiente es nuestra tabla de contingencia para las probabilidades conjuntas

```
##
##           Depression:0 Depression:1
## Misophonia:4  0.00000000  0.07317073
## Misophonia:3  0.20325203  0.04878049
## Misophonia:2  0.27642276  0.02439024
## Misophonia:1  0.04065041  0.00000000
## Misophonia:0  0.29268293  0.04065041
```

- ¿Cuál es la probabilidad marginal de misofonía de gravedad 3? (R/0.3)
- ¿Cuál es la probabilidad de no ser misofónico **y** no estar deprimido? (R/0.293)
- ¿Cuál es la probabilidad de ser misofónico **o** deprimido? (R/0.293)
- ¿Cuál es la probabilidad de ser misofónico **y** deprimido? (R/0.707)
- Describir en lenguaje escrito los resultados con probabilidad 0.

3.18.0.3 Ejercicio 3

Hemos realizado un experimento aleatorio 10 veces, que consiste en anotar el sexo y el estado vital de pacientes con algún tipo de cáncer después de 10 años del diagnóstico. Obtuvimos los siguientes resultados

```
##      A      B
## 1  male  dead
## 2  male  dead
## 3  male  dead
## 4 female alive
## 5  male  dead
## 6 female alive
## 7 female dead
## 8 female alive
## 9  male alive
## 10 male alive
```

- Crea la tabla de contingencia para el número $(n_{i,j})$ de observaciones de cada resultado (A,B)

- Crea la tabla de contingencia para la frecuencia relativa $(f_{i,j})$ de los resultados
- ¿Cuál es la frecuencia marginal de ser hombre? (R/0.6)
- ¿Cuál es la frecuencia marginal de estar vivo? (R/0.5)
- ¿Cuál es la frecuencia de estar vivo \cap ser mujer? (R/0.6)

3.18.0.4 Teoría: Ejercicio 4

- De la segunda forma de la regla de la suma, obtener la primera y la tercera forma.
- ¿Cuál es la regla de la suma de la tercera forma para la probabilidad de tres eventos $P(A \cup B \cup C)$?