

Master en Big Data. Fundamentos matemáticos del análisis de datos.

Sesión 3.

Fernando San Segundo

Curso 2020-21. Última actualización: 2020-09-01

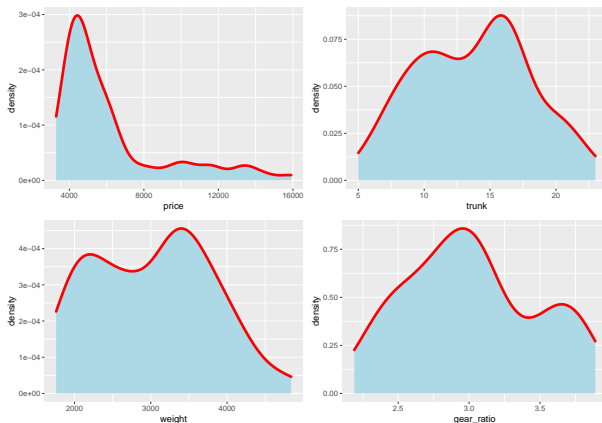


- 1 Distribuciones.
- 2 Valores centrales, de posición y dispersión.
- 3 Factores.
- 4 Comentarios finales.

Section 1

Distribuciones.

- En un ejercicio previo proponíamos dibujar las curvas de densidad de varias variables de la tabla `auto2`. Si lo haces obtendrás curvas como estas:

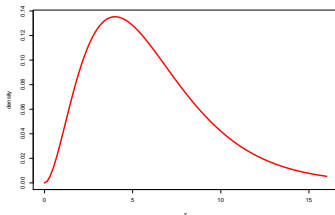


Mira el código de esta sesión para ver como dibujar esta “tabla de figuras”.

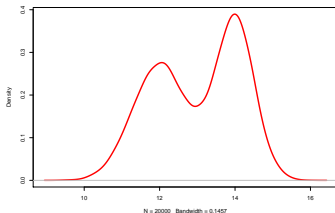
- Cada una de esas curvas muestra una *distribución de los valores*, que es una representación de la forma en la que se reparten los datos entre los valores posibles. Las distribuciones van a jugar un papel central en las próximas sesiones, así que vamos a desarrollar un poco de lenguaje para referirnos a ellas.

Distribuciones unimodales y multimodales.

- Una distribución que presenta un único máximo se denomina *unimodal*

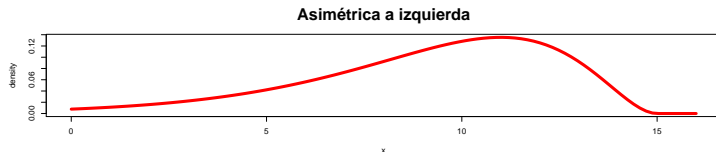
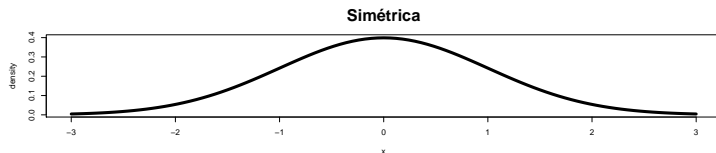
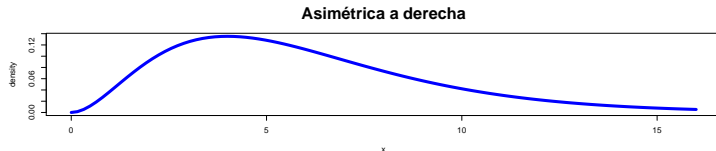


- Mientras que una distribución con dos máximos locales claramente definidos como la de la figura se denomina *bimodal* (o *multimodal* cuando son más de dos).



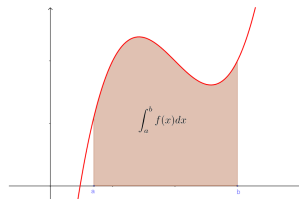
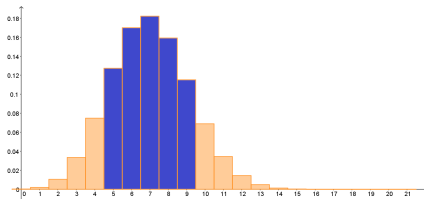
Distribuciones asimétricas.

- Otra característica de las distribuciones en la que nos vamos a fijar a menudo es su simetría. Fíjate en que es la *cola* más larga de la distribución la que da nombre a la asimetría. En inglés esta característica se denomina **skewness**



Distribuciones y cálculo de probabilidad. Discusión informal.

- La distribución de una variable discreta con valores enteros se representa adecuadamente mediante un histograma. Si las alturas de las barras del histograma representan frecuencias relativas entonces la probabilidad de que un valor elegido al azar caiga en el intervalo $[a, b]$ se obtiene sumando el área de todas las barras desde a hasta b , como se ilustra en la parte izquierda de la figura.
- Por su parte, La distribución de una variable continua se representa mediante una curva de densidad. Y entonces la probabilidad de que un valor elegido al azar caiga en el intervalo $[a, b]$ se obtiene calculando (con una integral) el área bajo la curva desde a hasta b



- Don't panic! No hace falta que sepas calcular integrales. Pero es importante entender que el cálculo de probabilidades está estrechamente relacionado con el cálculo de áreas.

Section 2

Valores centrales, de posición y dispersión.

La media aritmética

- Al trabajar con variables cuantitativas muchas veces trataremos de elegir un *valor central*. Es decir, un valor que sea *representativo* de los valores que toma la variable. Un buen valor central debería ser la respuesta a esta pregunta: “*si elijo un valor de la variable al azar ¿lo más probable es que se parezca a ... ?*”
- El valor central más conocido es la *media aritmética*. Dados n números x_1, x_2, \dots, x_n su media aritmética es:

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum_{i=1}^n x_i}{n}$$

- En R se utiliza la función `mean` para calcular la media de los valores de un vector numérico. Con este código vamos a elegir 20 números al azar entre 0 y 100 (con remplazamiento). Recuerda que los paréntesis sirven para que R muestre el resultado.

```
(muestra = sample(0:100, size = 20, replace = TRUE))
```

```
## [1] 24 41 68 16 60 75 18 25 8 89 45 71 22 93 71 66 83 46
```

```
## [19] 68 66
```

```
(media = mean(muestra))
```

```
## [1] 52.75
```

- Ejercicio:** Repítelo varias veces y mira las medias que obtienes. ¡Al pensar en esas medias estás empezando a hacer Estadística!

La media y los valores atípicos.

- La media aritmética se usa mucho porque su definición matemática es muy sencilla. Pero tiene un problema: su valor se ve muy afectado por la presencia de valores anormalmente grandes o pequeños de la variable, los llamados *valores atípicos* (que muy pronto vamos a definir con más rigor).
- Por ejemplo, si elegimos 99 números al azar de 0 a 100, su media es parecida a 50 como cabría esperar.

```
muestra = sample(0:100, size = 99, replace = TRUE)
(media = mean(muestra))
```

```
## [1] 51.64646
```

Pero si ahora añadimos un único valor igual a 1000 y volvemos a calcular la media:

```
muestra2 = c(muestra, 1000)
(media2 = mean(muestra2))
```

```
## [1] 61.13
```

Las dos muestras son esencialmente idénticas, así que al elegir un valor representativo (en el sentido probabilístico que hemos discutido) nos gustaría obtener respuestas mucho más parecidas. Pero como ilustra este ejemplo, la media se ve muy afectada por ese único valor atípico.

La mediana.

- Para paliar este problema podemos usar la *mediana* como sustituto o alternativa a la media. La mediana de un conjunto de valores se define con esta receta: se ordenan los números de menor a mayor y se toma el valor que queda en el centro (si hay una cantidad impar de valores; si son pares se promedian los dos valores centrales).
- Por ejemplo, para calcular la mediana de estos 17 valores:

```
(valores = sample(1:100, 17, replace = TRUE))
```

```
## [1] 25 42 69 17 61 76 19 26 9 90 46 72 23 94 72 67 84
```

los ordenamos y tomamos el valor central

```
(ordenados = sort(valores))  
(mediana = ordenados[9])
```

```
## [1] 9 17 19 23 25 26 42 46 61 67 69 72 72 76 84 90 94  
## [1] 61
```

Aunque lo mejor, claro, es usar la función de R:

```
median(valores)
```

```
## [1] 61
```

Más detalles sobre la mediana.

- Por su propia construcción (en términos de posiciones/rangos y no de tamaños) debería estar claro que la mediana no se ve muy afectada por la presencia de valores atípicos. Si volvemos a la muestra de 99 valores que construimos antes y miramos su mediana antes y después de añadir el valor atípico 1000 se obtiene:

```
median(muestra)
```

```
## [1] 54
```

```
median(muestra2)
```

```
## [1] 54
```

En este ejemplo concreto la mediana no se ve afectada en absoluto por ese valor.

- Y entonces, ¿por qué no se usa siempre la mediana en lugar de la media aritmética? Pues porque la definición de la mediana utiliza unas matemáticas bastante más complicadas que la media aritmética. Gracias al ordenador la importancia de los métodos basados en la mediana ha ido aumentando. Pero los métodos clásicos que usan la media aritmética siguen siendo los más comunes.
- En cualquier caso, recomendamos explorar bien los datos, identificar posibles valores atípicos y comprobar si la media y la mediana son parecidas.

- **Cuartiles y percentiles.** La mediana se puede pensar como el valor que es mayor o igual que el 50% de los valores de nuestro conjunto de datos. De la misma forma se definen el *primer y tercer cuartil* como los valores que son mayores o iguales respectivamente que el 25% y el 75% de los valores. Esto se generaliza a la noción de percentil. El percentil 43, por ejemplo, es un valor que deja por debajo al 43% de los valores. En R ese percentil se calcula para una variable como `cty` de `mpg` mediante:

```
quantile(mpg$cty, probs = 0.43)
```

```
## 43%
```

```
## 16
```

- El *recorrido* (del inglés *range*) de un conjunto de valores es el intervalo que va del valor mínimo al máximo. Cuidado: en español se usa a menudo erróneamente *rango*. Pero los rangos en Estadística son otra cosa (que corresponde a *rank* en inglés).
- Una manera sencilla de obtener varias de estas medidas de posición es usando `summary`

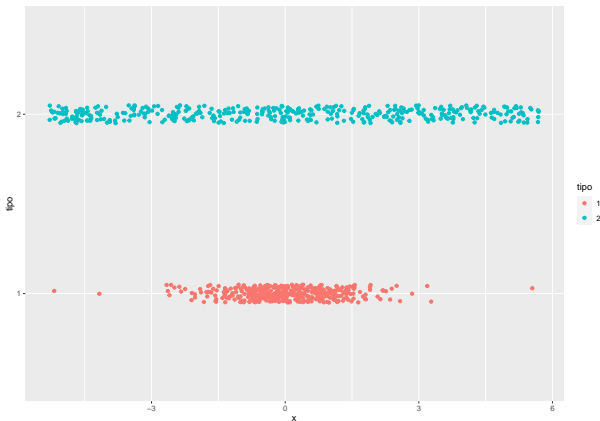
```
summary(mpg$cty)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      9.00  14.00   17.00   16.86   19.00   35.00
```

- **Ejercicio.** Prueba a aplicar la función a todo el `data.frame` con `summary(mpg)`.

Dispersión

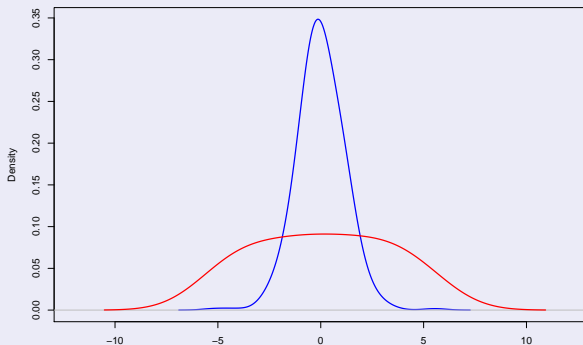
- La siguiente figura muestra dos conjuntos de valores, ambas con media 0 y el mismo número de puntos. Las coordenadas verticales de los puntos se han modificado con un poco de *ruido* para hacer mejor la visualización de ambos conjuntos. Puedes mirar el código de esta sesión para ver como se ha construido esta figura.



- ¿Qué diferencia a estos dos conjuntos de datos?

Dispersión en los diagramas de densidad.

- En esta otra figura puedes ver juntas las curvas de densidad de esos dos conjuntos de datos.



Este diagrama ilustra lo que la mostraban las nubes de puntos de la figura previa: los puntos de un conjunto están más agrupados en torno al centro que los del otro conjunto. La diferencia entre ambos conjuntos está en la **dispersión**.

- La idea intuitiva de la dispersión está clara: se trata de medir como de agrupados en torno al centro están los valores. ¿Como se define rigurosamente? Al igual que sucedía con los valores centrales (media vs mediana) hay varias posibilidades.

Recorrido intercuartílico. Valores atípicos.

- La primera forma de medir la dispersión que vamos a ver se basa en observar la diferencia entre los cuartiles tercero (75%) y primero (25%). Ese número es el *recorrido intercuartílico*; un intervalo de longitud IQR desde el primer hasta el tercer cuartil contiene siempre al 50% central de los valores. En R se calcula con:

```
IQR(mpg$cty)
```

```
## [1] 5
```

Fíjate en que es precisamente la diferencia entre esos dos cuartiles.

```
summary(mpg$cty)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      9.00  14.00   17.00   16.86   19.00   35.00
```

- El IQR nos permite dar una definición formal de los **valores atípicos**: son aquellos valores que superan $(\text{tercer cuartil}) + 1.5 \cdot \text{IQR}$ o quedan por debajo de $(\text{primer cuartil}) - 1.5 \cdot \text{IQR}$. Por ejemplo, en `mpg$cty` será atípico cualquier valor fuera de este intervalo

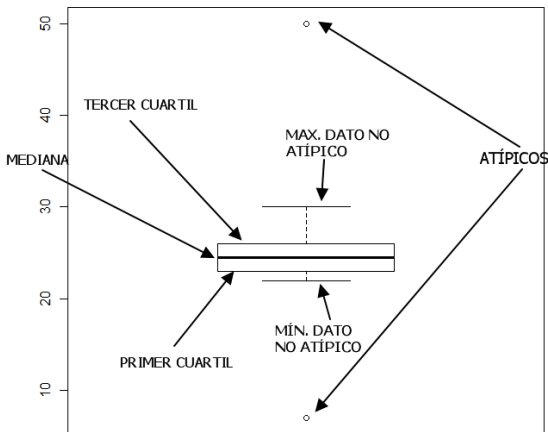
```
unnamed[quantile(mpg$cty, probs = c(1/4, 3/4)) + c(-1, 1) * 1.5 * IQR(mpg$cty)]
```

```
## [1] 6.5 26.5
```

- Ejercicio:** Mira el código de esta sesión y calcula el IQR de los dos conjuntos de datos de las páginas previas.

Boxplots

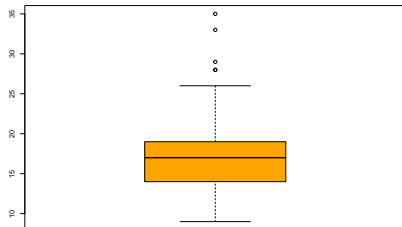
- El *boxplot* o diagrama de caja y bigotes (box and whiskers) es una forma de representar gráficamente estas medidas de posición. La estructura de un boxplot es la que se describe en esta figura:



Boxplots con R básico.

- Dibujar un boxplot con R básico es muy sencillo. Por ejemplo, para la variable `cty`:

```
bxp_cty = boxplot(mpg$cty, col="orange")
```



- Además al darle un nombre al boxplot podemos usar ese nombre para acceder a varios componentes del boxplot. Por ejemplo los valores atípicos en el vector de datos:

```
bxp_cty$out
```

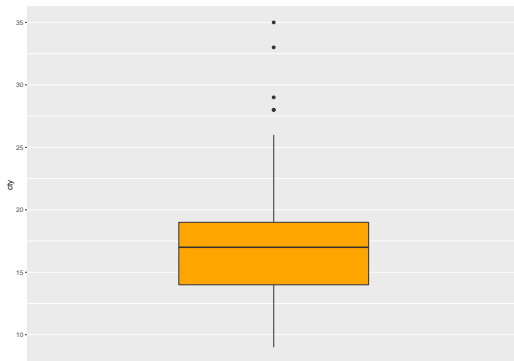
```
## [1] 28 28 33 35 29
```

- **Ejercicio:** ¿Cuáles son las *posiciones que ocupan* los valores atípicos de esta variable? Además haz un boxplot de `speed` en accidentes y busca en la ayuda del boxplot como dibujar el gráfico sin los valores atípicos.

Boxplot con ggplot

- En este caso usamos:

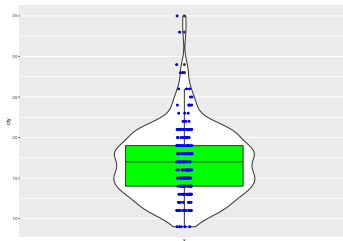
```
ggplot(mpg) +  
  geom_boxplot(mapping = aes(y = cty), fill="orange") +  
  scale_x_discrete(breaks = c())
```



Boxplot con datos y violinplots.

- Aunque los boxplots son muy útiles y se usan mucho, tienen sus limitaciones. La principal de ellas es que no muestran la distribución subyacente de valores (¡la *forma*!). Hay dos formas de combatir ese problema:
 - ▶ incorporando los valores al gráfico
 - ▶ usando un *violinplot* que es básicamente un boxplot con curvas de densidad añadidas.
- Esos dos remedios se pueden aplicar por separado o combinados como hemos hecho en este gráfico de ggplot:

```
ggplot(mpg) +  
  geom_violin(mapping = aes(x=0, y = cty)) +  
  scale_x_discrete(breaks = c()) +  
  geom_boxplot(mapping = aes(y = cty), fill="green") +  
  geom_jitter(aes(x=0, y = cty),  
    position = position_jitter(w=0.05, h= 0), col="blue")
```



Desviación absoluta mediana

- Además del rango intercuartílico se usan otras medidas de dispersión. Las más comunes están basadas en la idea de medir la *desviación absoluta* de cada dato individual con respecto a un valor central. Si los datos son x_1, x_2, \dots, x_n y el valor central (puede ser la media o la mediana) es c entonces las desviaciones absolutas son:

$$d_1 = |x_1 - c|, \quad d_2 = |x_2 - c|, \quad \dots \quad d_n = |x_n - c|$$

Usamos valores absolutos para evitar que las desviaciones por exceso compensen a las desviaciones por defecto. Aunque el valor absoluto es complicado...

- Para obtener una medida de dispersión buscamos un valor representativo (valor central) de las desviaciones absolutas. Por ejemplo, la mediana de las desviaciones absolutas respecto de la mediana (usando c igual a la mediana de los datos). Este valor es la **desviación absoluta mediana (MAD)**, que en R se calcula así:

```
library(readxl)
accidentes = read_excel("../datos/train_acc_2010.xls")
mad(accidentes$Speed, constant = 1)
```

```
## [1] 4
```

- Ejercicio:** calcula este valor a partir de los datos usando `median`. Mira lo que sucede si quitas la opción `constant = 1`. Lee la ayuda de `mad` y [esta página](#).
- Ejercicio:** calcula `mean(accidentes$Speed - mean(accidentes$Speed))`. Cambia de variable y de tabla y repite el cálculo. ¿Por qué pasa esto?

Varianza y desviación típica (definición poblacional).

- Usar los valores IQR o MAD como medidas de dispersión tiene inconvenientes similares a los que planteamos al comparar media y mediana. Las matemáticas de medianas y valores absolutos son complicadas y los resultados teóricos asociadas son más difíciles. Por esa razón tradicionalmente la Estadística ha usado una medida de dispersión basada en la media y con matemáticas más simples.
- Dados los valores numéricos x_1, x_2, \dots, x_n su **varianza (en sentido poblacional)** es la *media de las desviaciones cuadráticas respecto de la media aritmética \bar{x}* :

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

y la **desviación típica (también poblacional)** es la raíz cuadrada de la varianza:

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

Elevamos al cuadrado para conseguir un efecto similar al valor absoluto, pero usando una función derivable. Pero al hacerlo cambiamos las unidades y tenemos que tomar raíz cuadrada para obtener la dispersión en las unidades originales del problema.

Varianza y desviación típica muestrales.

- Al introducir la varianza y desviación típica les hemos puesto el apellido de *poblacionales*. Y tal vez te hayas dado cuenta de que no hemos dicho como calcular esos valores en R.
- La razón para hacer esto empezará a quedar más clara en las próximas sesiones y tiene que ver con la esencia misma de la Estadística. La misión fundamental de la Estadística se puede resumir en **obtener información fiable sobre la población a partir de una muestra**. En particular, para obtener información sobre la media y la dispersión en la población usaremos la media y la dispersión en la muestra. Y ese es el problema: usar en la muestra la misma fórmula de la varianza que en la población *no funciona bien*.
- El remedio es sencillo, afortunadamente. Hay que usar en la muestra las fórmulas *muestrales* para la varianza y desviación típica:

$$S_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1} \qquad s_x = \sqrt{s_x^2} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

como ves, la diferencia estriba en que dividimos por $n - 1$ en lugar de n .

Las funciones var y sd.

- Para calcular la varianza muestral en R usamos `var`:

```
var(mpg$displ)
```

```
## [1] 1.669158
```

puedes comprobar que esto es lo mismo que:

```
n = length(mpg$displ)
media = mean(mpg$displ)
sum((mpg$displ - media)^2) / (n - 1)
```

```
## [1] 1.669158
```

- Y para calcular la desviación típica muestral usamos `sd` (de *standard deviation*):

```
sd(mpg$displ)
```

```
## [1] 1.291959
```

puedes comprobar que esto es lo mismo que:

```
sqrt(var(mpg$displ))
```

```
## [1] 1.291959
```


Section 3

Factores.

Tablas de frecuencia y terminología para factores.

- Hasta ahora hemos centrado la discusión en variables cuantitativas, con valores numéricos. Vamos a tratar aquí muy brevemente las variables cualitativas o factores. Habrá muchas ocasiones en las próximas sesiones de practicar su uso.
- Al trabajar con factores sigue teniendo sentido usar tablas de frecuencias absolutas y relativas (pero no acumuladas, en general):

```
table(accidentes$TrkType)
```

```
##  
## Industry      Main Not rptd      Siding      Yard  
##          247      975          3          56      1340
```

```
prop.table(table(accidentes$TrkType))
```

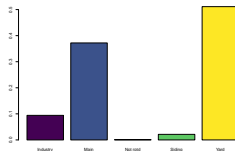
```
##  
##      Industry      Main      Not rptd      Siding      Yard  
## 0.094238840 0.371995422 0.001144601 0.021365891 0.511255246
```

- **¡Con factores genéricos los valores centrales o la dispersión no tienen sentido!**
- Los distintos valores de un factor se suelen llamar *niveles* (*levels*) del factor. Un factor *dicotómico* (o binario) tiene dos niveles; es decir, solo toma dos valores distintos. Si toma más de dos valores el factor es *politómico*. El nivel más frecuente es la *moda* del factor.

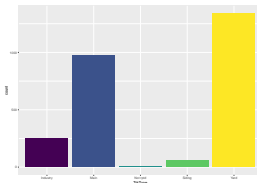
Gráficos para factores: diagramas de barras.

- El diagrama más simple y a menudo el más efectivo para representar la tabla de frecuencias de los niveles de un factor es el diagrama de barras. En R básico y ggplot respectivamente:

```
library(viridisLite)
barplot(prop.table(table(accidentes$TrkType)), col=viridis(5))
```



```
ggplot(accidentes) +  
  geom_bar(mapping = aes(x = TrkType), fill= viridis(5))
```



- Se desaconseja** el uso de diagramas de tarta.

Factores dicotómicos: variables binarias (de Bernoulli).

- Un factor dicotómico solo tiene dos niveles (valores). Esos valores se pueden representar como sí/no, cierto/falso, éxito/fallo, etc. Los factores dicotómicos también se llaman a menudo variables binarias o de Bernoulli.
- *Media como proporción.* Si pensamos en los niveles de un factor binario como cierto/falso, entonces al igual que hace R podemos identificar cierto/falso con 1/0. Esta representación numérica tiene una ventaja importante. Aunque hemos dicho que los valores centrales no tienen sentido en factores genéricos, en este caso concreto la media aritmética de un vector de unos y ceros *mide la proporción* de valores iguales a 1.
- Veamos un ejemplo. La variable `male` de la tabla `fhs` es un factor binario que vale 1 si el paciente es un hombre y 0 si es una mujer. Para averiguar la proporción de hombres en ese estudio basta con hacer:

```
fhs = read_csv("../datos/framingham.csv")  
mean(fhs$male)
```

```
## [1] 0.4292453
```

Aproximadamente el 43% de los pacientes son hombres.

Section 4

Comentarios finales.

Cadenas de caracteres (textos).

- Esta sesión ha tratado sobre los tipos de datos habituales en Estadística Clásica y que siguen siendo la base para la mayor parte de nuestro trabajo. Pero no queremos acabar sin mencionar que en los últimos años se ha producido una explosión en el tipo y volumen de datos que son objeto de análisis.
- Un ejemplo concreto de esos otros tipos de datos son los *textos* o cadenas de caracteres, que en R se corresponden con el tipo `character`. El análisis de datos de tipo texto juega un papel esencial en muchas aplicaciones como:
 - ▶ Buscadores
 - ▶ [Análisis de sentimiento](#)
 - ▶ [Sistemas de recomendación en comercio electrónico](#)
 - ▶ [Asistentes de voz](#) y muchas otras tecnologías actuales, conectadas con el Análisis del Lenguaje Natural (NLP) y la Inteligencia Artificial.
- R proporciona muchas herramientas para adentrarse en este tipo de problemas, e irás conociendo algunas en otras asignaturas. Aquí simplemente queremos mencionar que la librería [stringr](#), que es otro componente del *tidyverse*, es una de las mejores herramientas para iniciarse en el manejo de textos con R. Veremos ejemplos más adelante.

Rmarkdown para la creación de documentos.

- En las sesiones del curso veremos algunos ejemplos sencillos para que puedas iniciarte en el manejo de Rmarkdown.
- Rmarkdown, creado por Yihui Xie, es una herramienta muy general para la creación de documentos, en especial documentos relacionados con el Análisis de Datos. A partir de los ficheros escritos con Rmarkdown se pueden obtener fácilmente salidas en formato pdf o HTML, pero también presentaciones, artículos e informes técnicos, entradas de blog, libros, dashboards con Shiny y la lista sigue creciendo. Aprovecharemos las secciones finales de nuestras sesiones para presentar algunas de estas posibilidades.
- Hay dos fuentes de información online básicas para adentrarse en las posibilidades de RMarkdown:
Del propio Yihui Xie: [R Markdown: The Definitive Guide](#).
De RStudio: [R Markdown: Get Started](#).
- Recomendamos empezar leyendo una [introducción reciente a RMarkdon](#) de T. Goicoa (Univ. Pública de Navarra).

Enlaces

- [Código de esta sesión](#)
- [Cookbook for R](#)
- [Página web de ggplot2](#), que contiene el [resumen \(chuleta\)](#) elaborado por RStudio.
- [Resumen sobre importación de datos a R \(chuleta\)](#) elaborado por RStudio.
- Web del libro [PostData](#) y los tutoriales asociados. Para esta sesión se recomienda el Capítulo 2.

Bibliografía