

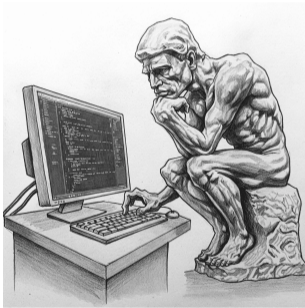
Análisis de datos en fortificación de alimentos a gran escala con R

Tema II: Importación de datos.

Dr. Maicel Monzón

Pregunta del encuentro anterior

Como se crea un objeto en R?



Respuesta

Los objetos se crean:

- Leyendo datos de un archivo
- Como resultado de un cálculo
- “Asignándoles un valor”
- etc.

Importar datos

Vamos a hablar de una de las formas más usadas para **crear objetos** en R.

Importar datos

Materia Prima de Datos



Los datos son esenciales para el análisis

Calidad de los Datos

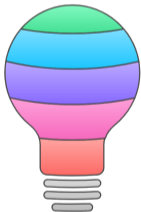


Validación inicial para asegurar la calidad

Colaboración



Facilita el trabajo en equipo y el acceso a los datos



Diversidad de Formatos

Los datos provienen de varios formatos



Automatización y Escalabilidad

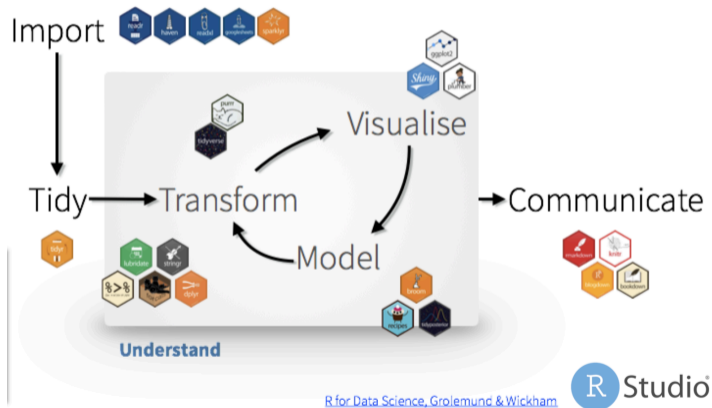
Procesos eficientes para proyectos grandes

Sumario

1. Introducción a la Importancia de Leer Datos
2. Tipos de Datos y Bibliotecas Especializadas
3. Funciones de readr para Lectura de Datos Tabulares
4. Convenciones y Argumentos Clave en readr
5. Parseo de Datos
6. Estrategias para Solucionar Problemas

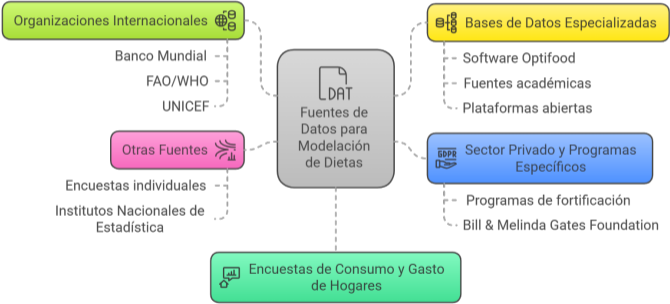
Lectura de datos en el proceso de ciencia de datos

Importar datos: Es fundamental porque los datos son la base de cualquier análisis.



Fuentes de datos para la modelación de dietas

Fuentes de Datos para Modelación de Dietas



Qué datos se pueden leer desde R?

- **Datos estructurados** (*estructura interna identificable filas, columnas con títulos*)

Ej. archivos rectangulares en texto plano (csv,dat,txt)

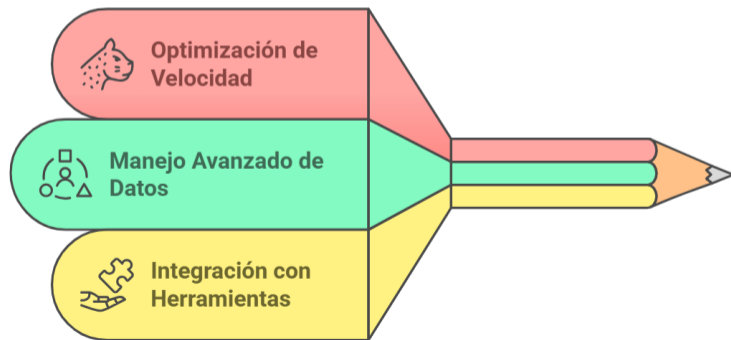
- **Datos no estructurados** (*datos binarios que no tienen estructura interna identificable*)

Ej. Correos electrónicos, publicaciones en redes sociales (Imágenes , sonido, video, etc.)

Bibliotecas especializadas en leer datos de distintos formatos

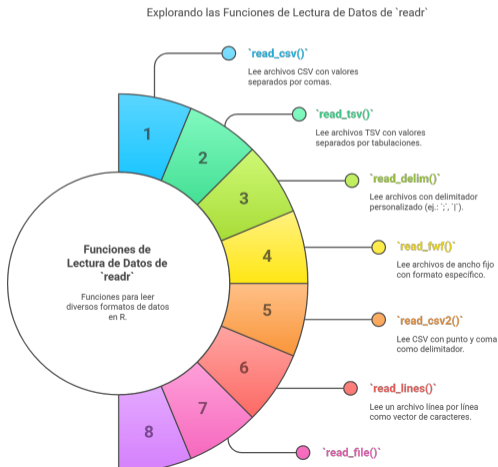
Paquete	formatos
readr*	datos rectangulares (csv, tsv y fwf)
haven	ficheros (SPSS, Stata y SAS)
readxl	ficheros Excel (.xls y .xlsx)
DBI	bases de datos
jsonlite	json
xml2	XML
httr	Web APIs
rvest	HTML (Web Scraping)

Bibliotecas como readr



Funciones para lectura de datos tabulares a tibbles (readr)

Estas funciones están diseñadas para leer datos rápidamente y devolver tibbles (data frames optimizados).



Delimitadores de readr

Funciones	delimitador
read_delim()	cualquier
read_csv()	coma
read_csv2()	punto y coma
read_tsv()	tabulaciones
read_fwf()	ancho fijo
fwf_widths()	ubicación

Algunas convenciones

path o camino (ruta hacia un archivo)

“./carpeta/archivo.ext”

separador entre carpetas es una barra inclinada (/),

(./) Espacio de trabajo

```
# obtiene el espacio de trabajo
ws <- getwd()
# contruye una ruta
file.path(".", "datos", "mifichero.csv")
```

Argumento file, readr(file)

```
read_csv(  
#file -cadena de texto- Argumento obligatorio  
file = "a,b,c \n  
      1,2,3 \n  
      4,5,6")
```

```
# A tibble: 2 x 3
```

	a	b	c
	<dbl>	<dbl>	<dbl>
1	1	2	3
2	4	5	6

Argumento skip, readr(file,skip)

```
read_csv(file = "La primer línea de metadata  
La segunda línea de metadata  
x,y,z  
1,2,3",  
#skip -omitir las primeras n líneas-  
skip = 2,  
col_types = "i" )
```

Argumento col_names

```
read_csv(file = "41,masculino, blanca  
         40, femenino, negra",  
# incluye el nombres de columna -tipo booleano-  
         col_names = FALSE)
```


Argumento col_name -Asignación de nombres de columna-

```
read_csv(file = "41,masculino,blanca\n
          40, femenino,negra",
# col_names -vector de caracteres-
          col_names = c("edad","sexo","color de la piel"))
```

Argumento na -valores faltantes-

```
# carácter que se asume como valores faltantes, vector de cadena  
read_csv(file = "a,b,c\n1,2, ", na = " ")
```

Parsing: proceso de convertir datos en formato correctos para su análisis

Parseo Transformación de Texto a Datos Estructurados



Lectura de Texto

Los datos se leen como cadenas de caracteres



Identificación de Tipo de Datos

Los datos se analizan para determinar el tipo



Conversión a Números

Las cadenas se convierten en valores numéricos



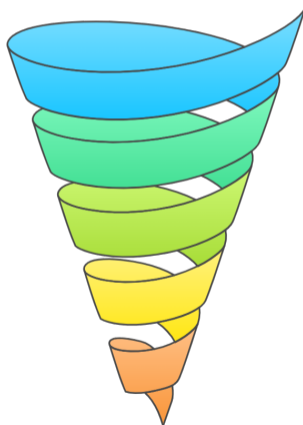
Conversión a Fechas

Las cadenas se convierten en formatos de fecha



Conversión a Booleanos

Las cadenas se convierten en valores verdadero/falso



funciones parse_*

funciones toman un vector de caracteres y devuelven:

- `parse_logical()` # **lógico**
- `parse_integer()` # **entero**
- `parse_double()` # **decimal**
- `parse_number()` # **numérico**
- `parse_character()` # **cadena**
- `parse_factor()` # **tipo factor**
- `parse_date()` # **fecha**
- `parse_time()` # **tiempo**

Parseo de números (Usos)

1. diferente de decimales *Ej coma o punto*

```
hemoglobina1<-c("12,2", "13.5", "11.9") # diferente marca decimal
```

2. números están rodeados por otros caracteres *Ej unidades*

```
hemoglobina2<-c("12.2", "13.5", "11.9 g/l") # texto adicionado
```

3. caracteres de “agrupación” *Ej “1,000,000”.*

Ejemplo 1: parse_double (locale) uso diferente de decimales

```
# locale -objeto que especifica las opciones de análisis que difieren de un
parse_double("1.23")
# decimal_mark "símbolo decimal"
parse_double("1,23",
             locale = locale(decimal_mark = ","))
```

Ejemplo 2: parse_number

```
# ignora los caracteres no-numéricos antes y después del número ignora los  
parse_number("123 g/l")  
parse_number("la hemoglobina es de 123 g/l")
```

Ejemplo 3 combinación de parse_number() y el locale

```
parse_number("123.456.789",  
             locale = locale(grouping_mark = "."))
```


Parseo de Cadenas de caracteres (parse_character())

```
x1 <- "El Ni\xf1o"  
# caracteres en espa\u00f1ol a veces no se leen bien ej \u00f1 o acentos  
parse_character(x1,  
                locale = locale(encoding = "Latin1"))
```

Parseo de factores parse_factor()

```
fruta <- c("manzana", "banana")  
# representar las variables categóricas que tienen un conjunto conocido de  
parse_factor(c("manzana", "banana", "bananana"), levels = fruta)
```

Parseo de fechas (parse_date())

- Año : %y (2 dígitos)
- Mes : %m (2 dígitos).
- Mes : %b (nombre abreviado, como “ene”).
- Mes :: %B (nombre completo, “enero”).
- Dias : “%d” (2 dígitos)

Parseo de tiempos (parse_time())

```
parse_time("20:10:01")
```

Analizar automatico de un archivo con readr (guess_parser)

```
# deduce automáticamente el formato con las primeras 1000 filas
guess_parser("2010-10-01")
guess_parser("15:01")
guess_parser(c("TRUE", "FALSE"))
guess_parser(c("1", "5", "9"))
guess_parser(c("12,352,561"))
guess_parser("12.3")
```

Detección automática (1000 primeras filas)

`guess_parser` casi siempre detecta los tipos correctos, sin embargo, puede no ser así en todos los casos.




función `problems()`: Identificación de problemas de parseo

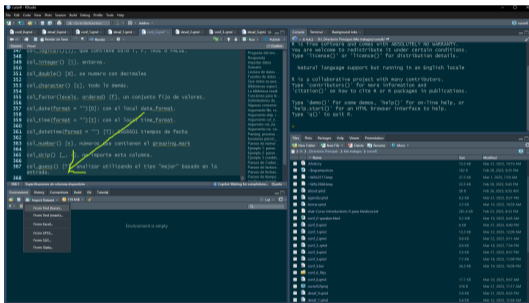
```
desafio <- read_csv("datos_fortificacion.csv")  
# si sale notificaciones  
problems(desafio)
```

Funciones para Definir Tipos de Columnas en readr

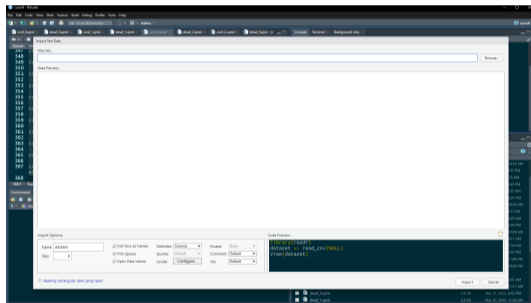
Cada función `col_*()` se usa para indicar el tipo de dato de una columna al importar datos. Las **abreviaturas** (entre paréntesis) permiten definir tipos de forma rápida.

Función	Abreviatura	Descripción
<code>col_logical()</code>	[l]	Valores lógicos: TRUE, FALSE, T, F.
<code>col_integer()</code>	[i]	Números enteros (ej: -5, 8, 100).
<code>col_double()</code>	[d]	Números con decimales (ej: 3.14, -9.5).
<code>col_character()</code>	[c]	Texto (se usa si hay caracteres no numéricos). Valor por defecto en readr.
<code>col_factor(levels, ordered)</code>	[f]	Factores: categorías con niveles definidos (ej: "hombre", "mujer").
<code>col_date(format = "")</code>	[D]	Fechas (ej: "2023-12-25"). Usa el formato de fecha del <code>locale</code> .
<code>col_time(format = "")</code>	[T]	Horas (ej: "14:30:00"). Usa el formato de hora del <code>locale</code> .
<code>col_datetime(format = "")</code>	[DT]	Fechas y horas (ej: "2023-12-25 14:30:00"). Formato ISO8601 por defecto.
<code>col_number()</code>	[n]	Números con separadores de miles (ej: "1,000.50" → 1000.5).
<code>col_skip()</code>	[.] o [-]	Omite la columna (no la importa).
<code>col_guess()</code>	[?]	Tipo automático (lo que hace readr por defecto si no es específica). 

Importar con interfaz gráfica



Interfaz gráfica para la función readr



Conclusión

La **importación de datos** es el **primer paso crítico** en cualquier **análisis**. Con **readr** y **bibliotecas especializadas**, R ofrece herramientas eficientes para manejar **múltiples formatos**. Dominar estos recursos permite **enfocarse en el análisis** y la **toma de decisiones basada en datos**.

¡Gracias por su atención!