

# Fundamentos de Análisis de Datos

## **Unidad 3**

Extracción de datos

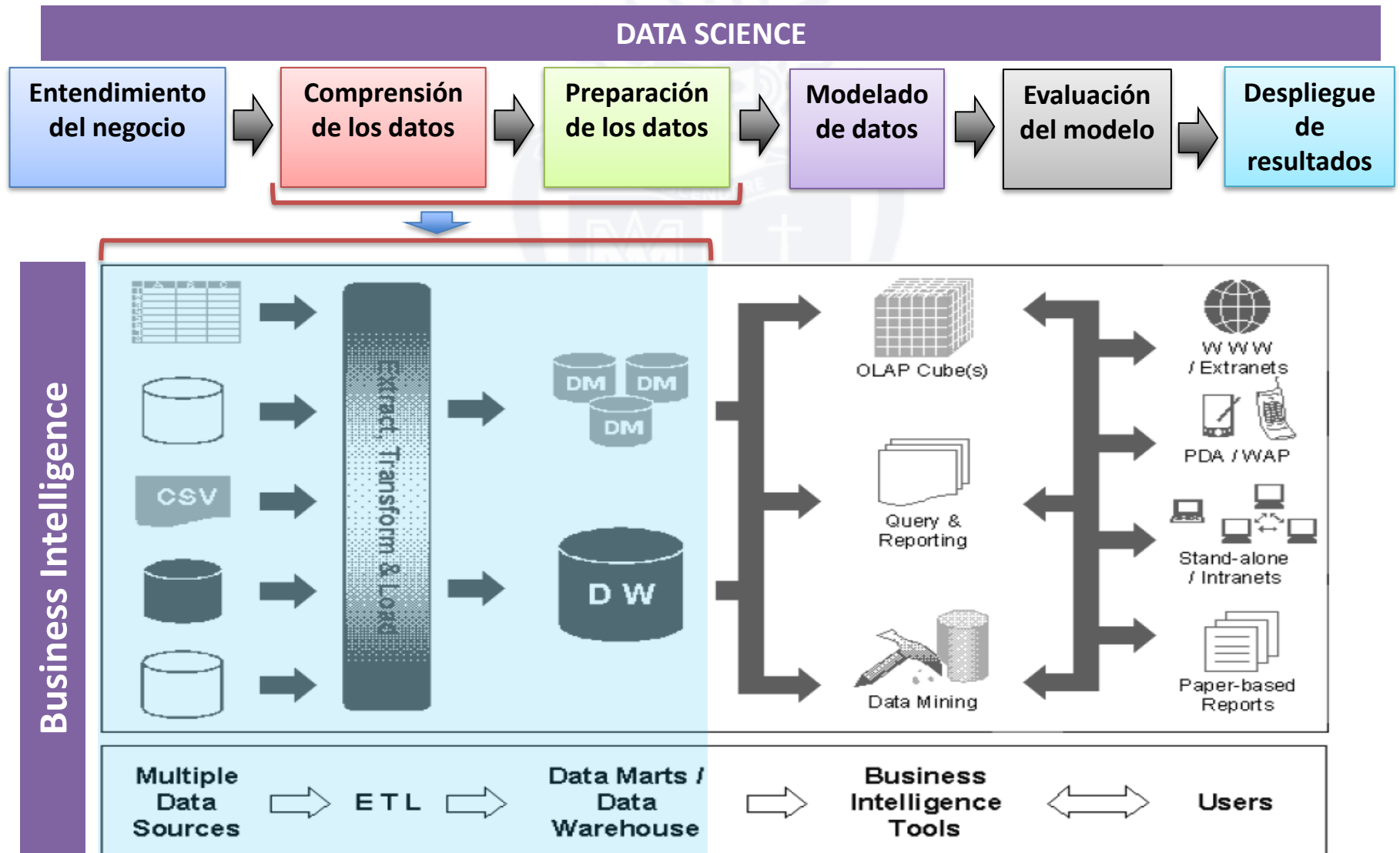
UTPL

UNIVERSIDAD TÉCNICA PARTICULAR DE LOJA

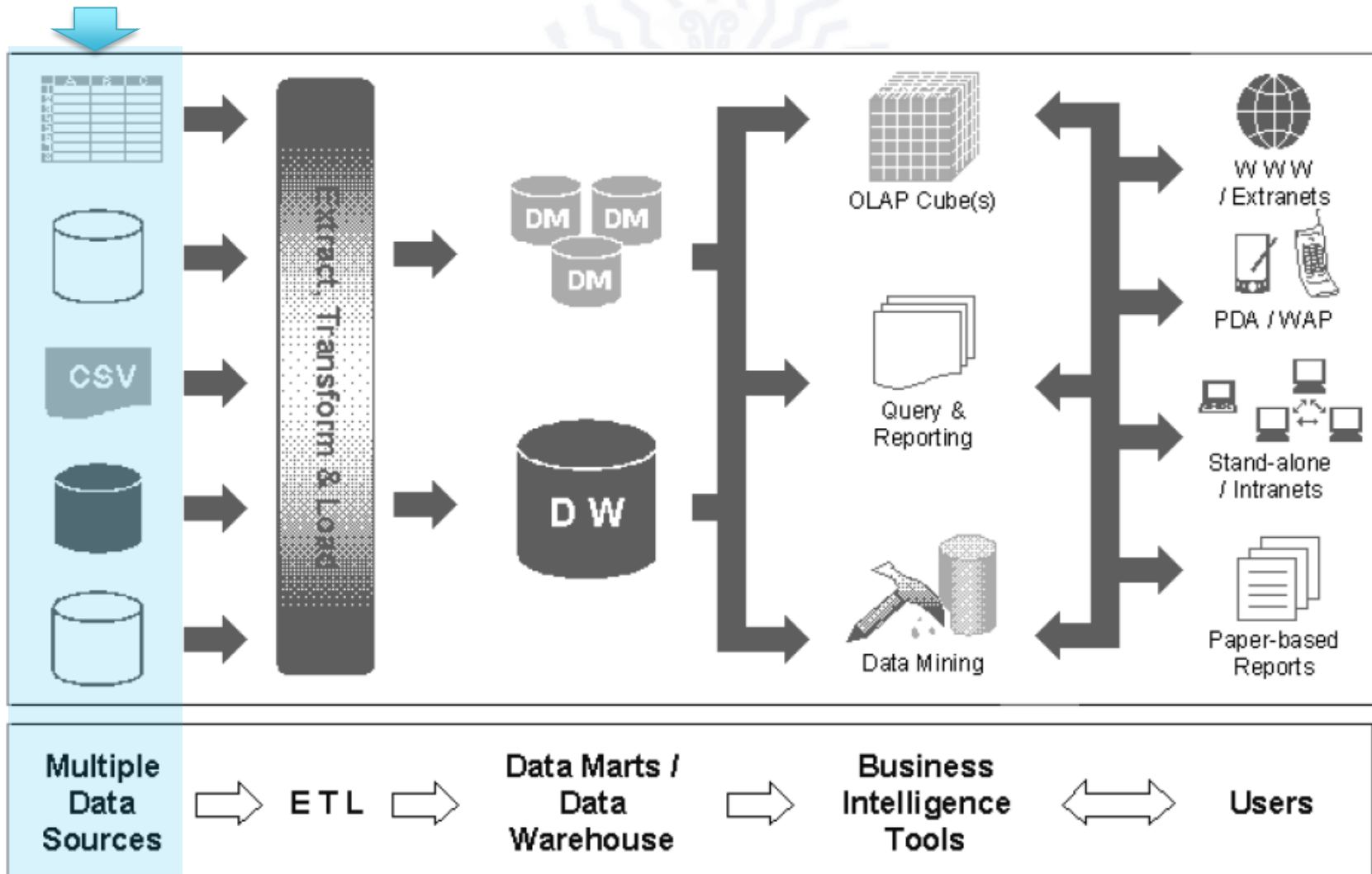
# Agenda

- El proceso ETL
- Librerías para extracción y manipulación de datos
- Librerías para el manejo de datos
- Lectura, consumo y visualización de datos
- Técnicas de extracción de datos
  - Extracción desde archivos estructurados

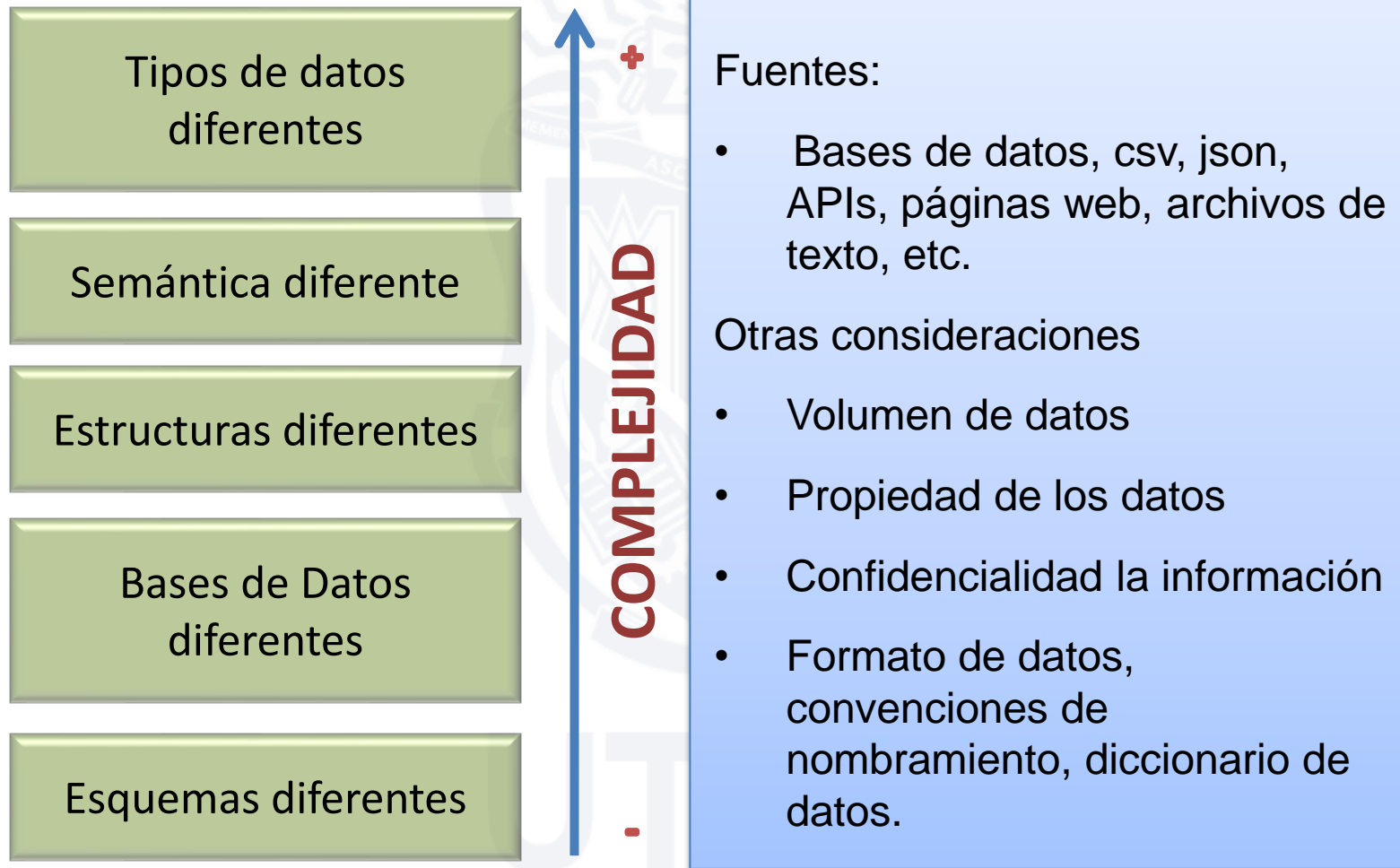
# El proceso ETL



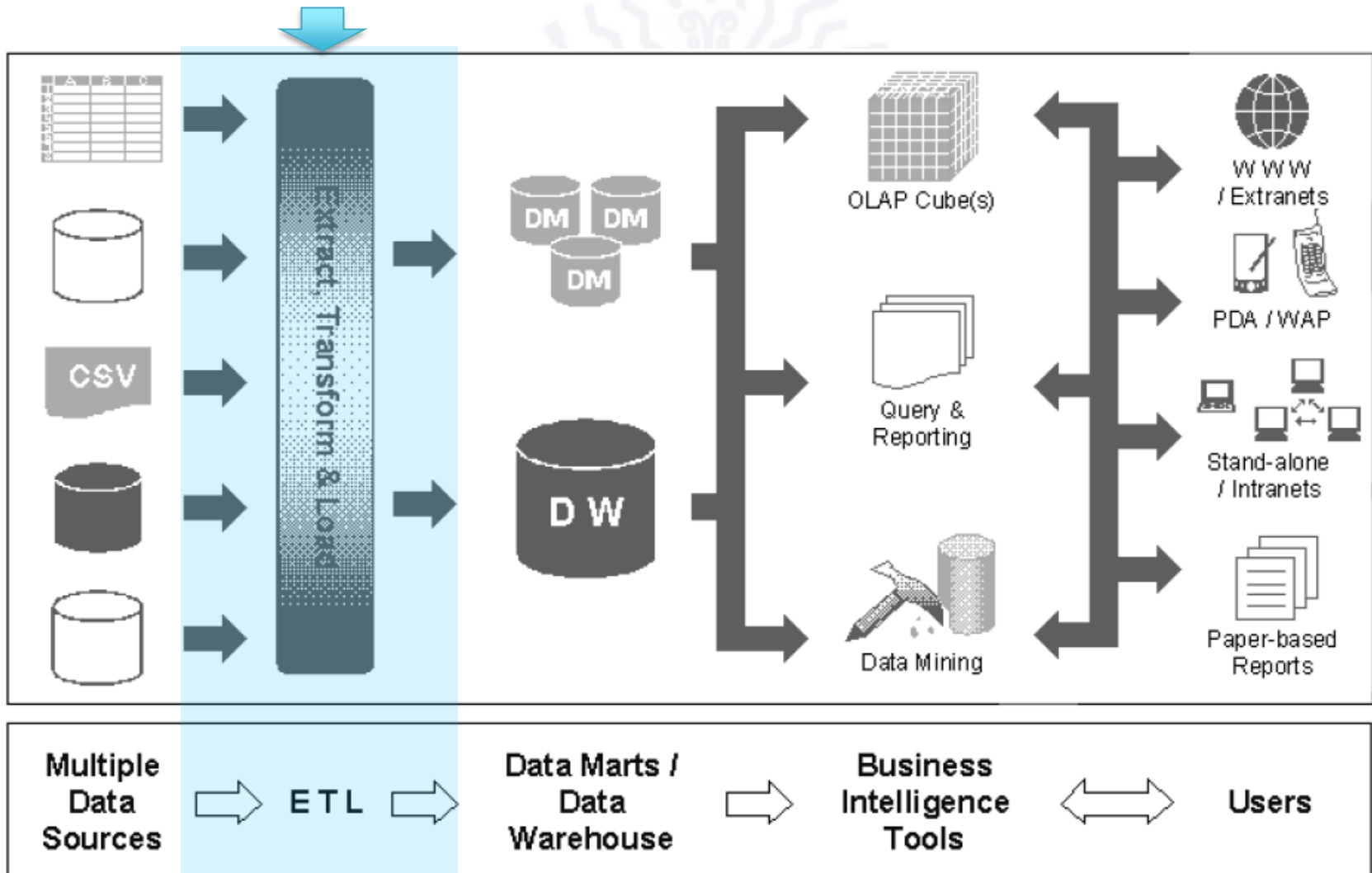
# Fuentes de datos



# Fuentes de Datos



# Extracción, Transformación y Carga (ETL)



# Extracción, Transformación y Carga (ETL)

## EXTRACCIÓN

- Estudio de los datos: Permitirá identificar:
  - Datos fuera de los límites (atípicos).
  - Violaciones a las dependencias (Ej campos derivados).
  - Datos redundantes.
  - Datos huérfanos (integridad referencial).
- Extracción de los datos:
  - En qué formatos se encuentran los datos?
  - Cuál será la frecuencia con la que se extraerán los datos?
  - En qué orden se cargarán los datos?
  - Cómo se minimizará el tiempo requerido para cargar los datos?
- Estacionamiento de los datos
  - Almacén de datos temporal
  - Evitar sobrecarga en las fuentes de datos

Extracción  
eficiente

Replica fiel  
de datos  
origen

Evitar  
impacto en  
las fuentes

# Extracción, Transformación y Carga (ETL)

## TRANSFORMACIÓN

Convertir los  
datos brutos  
en un formato  
más adecuado  
para análisis

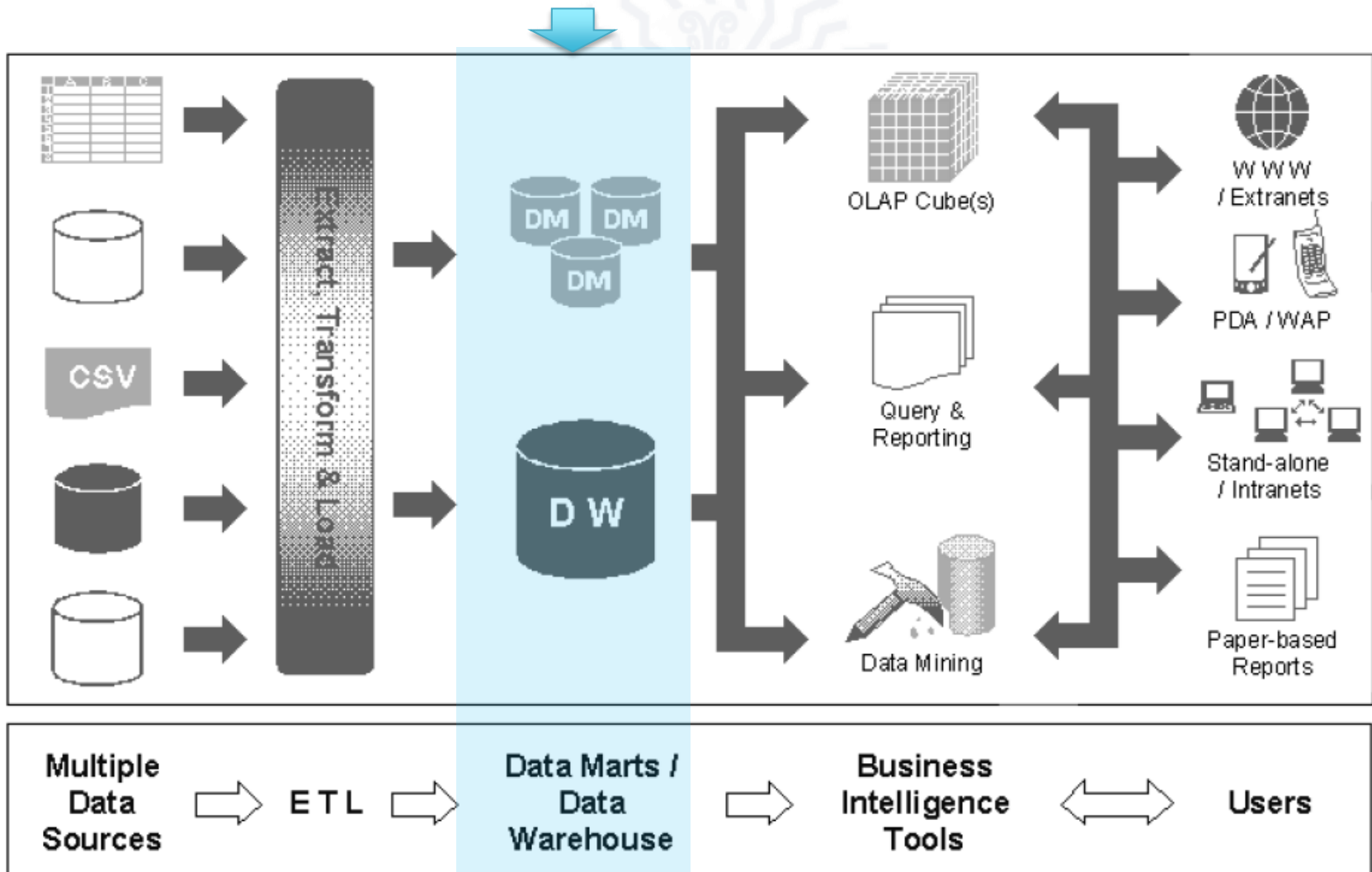
Datos  
homogéneos  
libres de  
inconsistencias

**CALIDAD**

- Limpieza de los datos
  - Ausencia de valores
  - Valores inconsistentes
  - Valores duplicados
  - Violaciones a las reglas del negocio
- Enriquecimiento de datos
  - Ordenar
  - Dividir
  - Juntar
  - Derivar
  - Buscar y Anexar
- Agrupación de datos (resúmenes)



# Almacenamiento integrado de datos



# Extracción, Transformación y Carga (ETL)

## CARGA



# Importancia del proceso ETL

- **Facilita la integración de los datos:** Consolidación de datos de múltiples fuentes.
- **Mejora la calidad de los datos:** Limpia, corrige, estandariza, enriquece los datos
- **Soporte de Decisiones:** Alimentan sistemas de inteligencia de negocios y análisis.
- **Mejora del Acceso a los Datos:** Acceso más rápido y eficiente a los datos unificados y organizados.
- **Escalabilidad:** Puede escalar para manejar volúmenes de datos mayores.
- **Eficiencia Operativa:** Automatiza el proceso de mover y transformar datos, reduciendo errores.
- **Cumplimiento y Seguridad de Datos:** Ayuda al cumplimiento de normativas y políticas de seguridad en el tratamiento de datos.

# Librerías para extracción y manipulación de datos

- R:
  - **xlsx**: manejo de archivos excel
  - **DBI**: para interactuar con bases de datos en R, lo que facilita la conexión y la ejecución de consultas
  - **httr**: para la extracción de datos de APIs y otros recursos en línea
  - **rvest**: para web scraping
  - **dplyr**: facilita la selección, filtrado de filas, reordenamiento y agregación de datos
  - **tidyr**: permite transformar los datos en un formato más adecuado para análisis. Reorganizar, completar, separar, unir columnas
  - **jsonlite**: Leer y escribir archivos JSON.
  - **lubridate**: Manejar fechas y tiempos.
  - **janitor**: Limpiar nombres de columnas, detectar duplicados y valores inconsistentes.
  - **validate**: Validar reglas de calidad de datos.
  - **stringr**: Operaciones sobre texto (buscar, reemplazar, extraer).

# Librerías para extracción y manipulación de datos

- Python:
  - **Pandas**: Es esencial para la manipulación de los datos: limpieza, transformación y análisis
  - **Numpy**: permite manipular estructuras de datos de tipo matriz
  - **requests**: Consumir APIs y manejar peticiones HTTP.
  - **beautifulsoup4**: Web scraping de HTML/XML.
  - **Scrapy**: Framework de web scraping a gran escala.
  - **SQLAlchemy**: ORM para conectar Python con bases de datos relacionales.
  - **openpyxl**: Trabajar con archivos Excel (.xlsx).
  - **json**: Trabajar con JSON nativo en Python.
  - **Selenium**: automatizar la interacción con navegadores web. Web scraping avanzado

# Lectura, consumo y visualización de datos

## **Lectura de datos:**

Es el proceso de importar datos almacenados en archivos o bases de datos hacia un entorno de análisis, para trabajar con ellos de manera local.

## **Consumo de datos:**

Es la obtención de datos provenientes de fuentes dinámicas, como APIs o servicios en línea, accediendo a información actualizada bajo demanda o en tiempo real.

## **Visualización de datos:**

Es la transformación de datos en representaciones gráficas, facilitando la identificación de patrones, tendencias y la comunicación efectiva de resultados.

# Lectura, consumo y visualización de datos

Concepto	Lectura de datos	Consumo de datos
Qué es	Es el acto técnico de importar datos desde un origen a tu programa o entorno de trabajo.	Es el proceso de adquirir, procesar y actualizar datos, usualmente de fuentes dinámicas o remotas.
Enfoque	Recuperar un archivo o dataset completo para comenzar a trabajar.	Obtener datos que pueden cambiar frecuentemente o que requieren consultas activas.
Ejemplos típicos	Leer archivos CSV, Excel, JSON locales. Leer tablas completas de bases de datos.	Consultar APIs con parámetros. Hacer web scraping de páginas en tiempo real. Actualizar datos periódicamente desde bases de datos o servicios.
Frecuencia	Generalmente puntual: una vez o al iniciar un proyecto.	Puede ser continua, en intervalos o en tiempo real.
Tipo de fuente	Normalmente estática.	Normalmente dinámica.
Objetivo	Traer el conjunto de datos para su posterior análisis.	Integrar datos como parte activa de un flujo de trabajo o aplicación.

# Lectura, consumo y visualización de datos

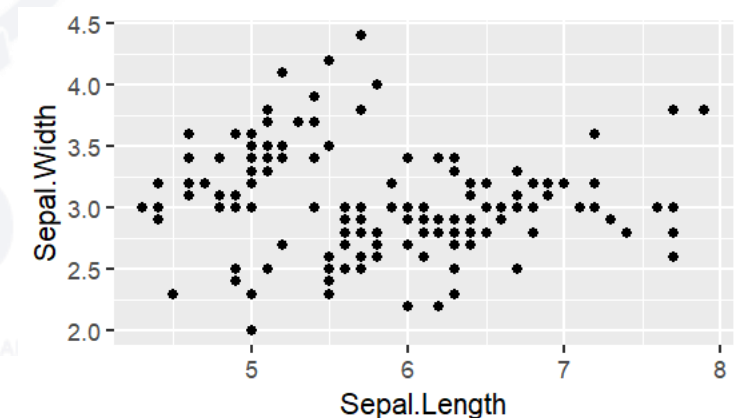
La **visualización de datos** convierte la información en representaciones gráficas para facilitar su comprensión y análisis.

Tipos de visualizaciones comunes:

- **Gráficos de barras:** comparar cantidades.
- **Histogramas:** analizar distribución de datos.
- **Diagramas de dispersión:** identificar correlaciones entre variables.
- **Mapas de calor (heatmaps):** detectar concentraciones de valores.
- **Gráficos de líneas:** mostrar evolución a lo largo del tiempo.

Herramientas principales (R):

- **ggplot2:** gráficos estáticos de alta calidad.
- **plotly:** gráficos interactivos.
- **leaflet:** mapas interactivos.





# Técnicas de extracción de datos

- Extracción de datos desde archivos estructurados
- Extracción de datos desde bases de datos
- Extracción de datos desde archivos semi-estructurados
- Extracción de datos mediante web scraping
- Extracción de datos usando APIs
- Extracción de datos desde archivos SAS

# Extracción desde archivos estructurados

- Se pueden extraer datos desde archivos estructurados (ejemplo: csv, excel)
- En R se lo puede hacer mediante funciones que trae el paquete **base**, y mediante otras librerías como **xlsx**.
- Es una de las formas más comunes de compartir y adquirir datos.

# Extracción desde archivos estructurados

## Ejemplo de extracción de datos desde archivos CSV

```
3  
4 # Configurar el directorio de trabajo  
5 setwd("C:/Users/aeenc/RStudio")  
6  
7 # Leer y cargar el contenido del archivo csv  
8 datos_csv <- read.csv("titanic.csv", sep = ",", quote = "\"", dec = ".")  
9  
10 # Ver el contenido del data frame  
11 print(datos_csv)  
12 str(datos_csv)
```



# Extracción desde archivos estructurados

## Ejemplo de extracción de datos desde archivos EXCEL

```
4 # Configurar el directorio de trabajo
5 setwd("C:/Users/aeenc/RStudio")
6
7 # Instalar y cargar librería "xlsx"
8
9 install.packages("xlsx")
10 library(xlsx)
11
12 # Leer y cargar el contenido de la primera hoja del archivo xlsx
13 datos_excel <- read.xlsx("titanic.xlsx", sheetIndex = 1)
14
15 # Ver el contenido del data frame
16 print(datos_excel)
17 str(datos_excel)
18
```

# dplyr

**dplyr** es una librería de R que facilita la **manipulación de datos**. Permite **seleccionar, filtrar, ordenar, agrupar y resumir** datos de manera rápida y legible.

Se puede usar con el operador **%>%** (pipe) para conectar pasos de manipulación de **izquierda a derecha**

## Funciones básicas de dplyr

Función	¿Qué hace?
filter()	Filtra filas según condiciones
select()	Selecciona columnas específicas
arrange()	Ordena filas
mutate()	Crea o transforma columnas
summarise()	Resume datos (agrega, promedia, etc.)
group_by()	Agrupar datos para operaciones por grupo

# Extracción desde bases de datos

- Conlleva la adquisición de datos desde bases de datos relacionales o no relacionales.
- Requiere establecer una conexión al servicio de base de datos, y el envío de una consulta mediante lenguaje de base de datos (Ej: SQL)
- En R se puede usar la librería **DBI**.
- **DBI** permite trabajar con varios tipos de bases de datos. Puede requerir instalar paquetes adicionales de acuerdo al motor específico.

# Extracción desde bases de datos

## Ejemplo de extracción de datos desde MySQL

```
2 # Instalación y carga de librerías
3 install.packages("DBI")
4 install.packages("RMySQL")
5
6 library("DBI")
7 library("RMySQL")
8
9 # Conexión a una base de datos MySQL
10 con = dbConnect(RMySQL::MySQL(),
11                 host = "localhost",
12                 dbname = "enemdu",
13                 user = "user01",
14                 password = "1q2w3e4r",
15                 port = 3306)
16
17 # Construcción de la consulta SQL
18 comando_sql = paste0("select * ",
19                       "from encuestas ",
20                       "where cod_provincia = 11 ",
21                       "and genero = \"H\"; ")
22
23 # Ejecución de la consulta SQL
24 datos_mysql <- dbGetQuery(con,comando_sql)
25
26 #Cierre de la conexión MySQL
27 dbDisconnect(con)
28
29 View(datos_mysql)
```

# Extracción desde bases de datos

## Ejemplo de extracción de datos desde Oracle

```
2 # Instalación y carga de librerías
3 install.packages("odbc")
4
5 library("odbc")
6
7 # Búsqueda de un DSN asociado al DBMS a conectar
8 odbc::odbcListDrivers()
9
10 # Conexión a Oracle via ODBC
11 con <- dbConnect(odbc::odbc(),
12                 Driver = "Oracle in XE",
13                 Server = "localhost",
14                 UID = "hr",
15                 PWD = "oracle",
16                 Port = 1521,
17                 SVC = "XE")
18
19 # Construcción de la consulta SQL
20 comando_sql = "select * from employees;"
21
22 # Ejecución de la consulta SQL
23 datos_oracle = dbGetQuery(con,comando_sql)
24
25 #Cierre de la conexión MySQL
26 dbDisconnect(con)
27
28 View(datos_oracle)
```



# Extracción desde archivos semi-estructurados

- Se pueden también extraer datos desde archivos semiestructurados (Ejemplo: JSON)
- En R se lo puede hacer mediante librerías como **jsonlite**.
- Es otra de las formas más comunes de compartir y adquirir datos.
- JSON se usa también en el contexto de extracción de datos a partir de APIs.

# Extracción desde archivos semi-estructurados

## Ejemplo 1 de extracción de datos desde archivos JSON

```
2 # Instalación y carga de librerías
3 install.packages("jsonlite")
4
5 library(jsonlite)
6
7 # Leer y cargar contenido desde archivo json
8 libros <- fromJSON("LibrosPearson.json")
9
10 # Inspeccionar formato y contenido de dataset resultante
11 class(libros)
12 View(libros)
13 str(libros)
14 libros$area
15 libros$autor
```

# Extracción desde archivos semi-estructurados

## Ejemplo 1 de extracción de datos desde archivos JSON

```
1 [
2   {
3     "_id": "9786073233316",
4     "titulo": "CALCULO UNA VARIABLE",
5     "edicion": 13,
6     "area": "Matemáticas",
7     "año": 2015,
8     "autor": "THOMAS"
9   }, {
10    "_id": "9786073235761",
11    "titulo": "GEOMETRIA ANALITICA",
12    "edicion": 4,
13    "area": "Matemáticas",
14    "año": 2015,
15    "autor": "AGUILAR"
16  }, {
17    "_id": "9786073235808",
18    "titulo": "GEOMETRIA, TRIGONON",
19    "edicion": 4,
20    "area": "Matemáticas",
21    "año": 2015,
22    "autor": "AGUILAR"
23  }, {
24    "_id": "9786073235853",
25    "titulo": "CALCULO DIFERENCIAL",
26    "edicion": 4,
27    "area": "Matemáticas",
28    "año": 2015, "autor": "AGUILAR"
```

*LibrosPearson.json*



*Data frame LIBROS*

	_id	titulo	edicion	area	año	autor
1	9786073233316	CALCULO UNA VARIABLE	13	Matemáticas	2015	THOMAS
2	9786073235761	GEOMETRIA ANALITICA	4	Matemáticas	2015	AGUILAR
3	9786073235808	GEOMETRIA, TRIGONOMETRIA Y GEOMETRIA ANALITICA	4	Matemáticas	2015	AGUILAR
4	9786073235853	CALCULO DIFERENCIAL E INTEGRAL	4	Matemáticas	2015	AGUILAR
5	9786073237451	ALGEBRA LINEAL Y SUS APLICACIONES	5	Matemáticas	2016	C. LAY
6	9786073237642	MATEMATICAS 1	1	Matemáticas	2016	c("ESTRADA", "JIMENEZ")
7	9786073237789	TALLER DE LECTURA Y REDACCION 1	3	Humanidades	2016	DE TERESA
8	9786073238021	COMO PROGRAMAR EN JAVA	10	Computación	2016	DEITEL
9	9786073238229	FISICA CONCEPTUAL	12	Ciencias	2016	HEWITT
10	9786073238380	BIOLOGIA	1	Ciencias	2016	GAMA
11	9786073239073	INFORMATICA 2	1	Computación	2016	PÉREZ
12	9786073239233	QUIMICA 2	1	Ciencias	2016	GUTIERREZ
13	9786073240994	MECANICA DE MATERIALES	9	Ingeniería	2017	HIBBELER
14	9786073241984	ORIENTACION EDUCATIVA IV	2	Humanidades	2018	ONOFRE
15	9786073242011	EL PLACER DE LA ESCRITURA MANUAL DE APROPIACION D...	5	Humanidades	2018	CORREA

# Extracción desde archivos semi-estructurados

## Ejemplo 2 de extracción de datos desde archivos JSON

```
4  
5 library(jsonlite)  
6  
7 # Leer y cargar contenido desde archivo json  
8 personas <- fromJSON("personas.json",flatten = TRUE)  
9  
10 # Inspeccionar formato y contenido de dataset resultante  
11 class(personas)  
12 view(personas)  
13 str(personas)
```



# Extracción desde archivos semi-estructurados

## Ejemplo 2 de extracción de datos desde archivos JSON

```
1 [{
2   "cedula": "111111111",
3   "nombre": "juan",
4   "telefono": [
5     "0978786767",
6     "0728938293"
7   ],
8   "sueldo": 1200
9 }, {
10  "cedula": "222222222",
11  "nombre": "maria",
12  "telefono": [
13    "027428483",
14    "0983764552"
15  ],
16  "sueldo": 1230.34,
17  "titulo": {
18    "denominacion": "Arquitecto",
19    "anio": 1990
20  }
21 }, {
22  "cedula": "333333333",
23  "nombre": "pedro",
24  "telefono": "0498263541",
25  "titulo": {
26    "denominacion": "Ingeniero en informática",
27    "anio": 2000
28  }
29 },
```



*Data frame PERSONAS*

	cedula	nombre	telefono	sueldo	email	titulo.denominacion	titulo.anio
1	111111111	juan	c("0978786767", "0728938293")	1200.00	NA	NA	NA
2	222222222	maria	c("027428483 ", "0983764552")	1230.34	NA	Arquitecto	1990
3	333333333	pedro	0498263541	NA	pedro@mail.com	Ingeniero en informática	2000

*Personas.json*

# Extracción mediante web scraping

- Es posible extraer datos desde una página web mediante la extracción de elementos de la página, con base en la lectura de nodos dentro del código HTML.
- Se usa cuando no podemos obtener los datos mediante una descarga en formatos convencionales (CSV, JSON, etc.), o mediante APIs.
- En R para este proceso se pueden usar funciones que vienen con la librería “**rvest**”
- Web scraping debe realizarse considerando las leyes de privacidad y derechos de autor, así como los términos de servicio del sitio web.
- Debe realizarse de manera responsable evitando impacto en las fuentes

# Extracción mediante web scraping

## Principales funciones **rvest**

Función	Descripción breve
<code>read_html()</code>	Lee el HTML de una página web o archivo local
<code>html_elements()</code>	Extrae múltiples nodos HTML que coinciden con un <b>selector</b>
<code>html_element()</code>	Extrae el primer nodo que coincide con un <b>selector</b>
<code>html_text()</code> / <code>html_text2()</code>	Extrae el texto interno del nodo HTML
<code>html_attr()</code>	Extrae el valor de un atributo HTML (ej. href, src, title)
<code>html_table()</code>	Convierte tablas HTML en data.frames

UNIVERSIDAD TÉCNICA PARTICULAR DE LOJA

# Extracción mediante web scraping

## ¿Qué es un selector CSS?

- Un **selector CSS** es una regla que permite identificar elementos dentro del código HTML de una página web, basándose en atributos como **etiqueta (tipo)**, **clase**, **ID**, o **posiciones estructurales**.
- En el contexto de **rvest en R**, un selector CSS es una **cadena de texto** que le indica al programa **qué elementos del dominio queremos extraer**. Funciona como una **guía precisa** para señalar contenido como títulos, listas de precios, enlaces, imágenes, entre otros.
- Es equivalente a decirle a R:
  - “Encuentra todos los elementos <h1>”,
  - “Extrae todas las imágenes con clase producto”, o
  - “Tráeme los enlaces dentro de una tabla”.
- Un selector CSS en rvest actúa como una **fórmula estructurada para localizar elementos HTML**, y se utiliza en funciones como `html_nodes()` o `html_elements()`.



# Extracción mediante web scraping

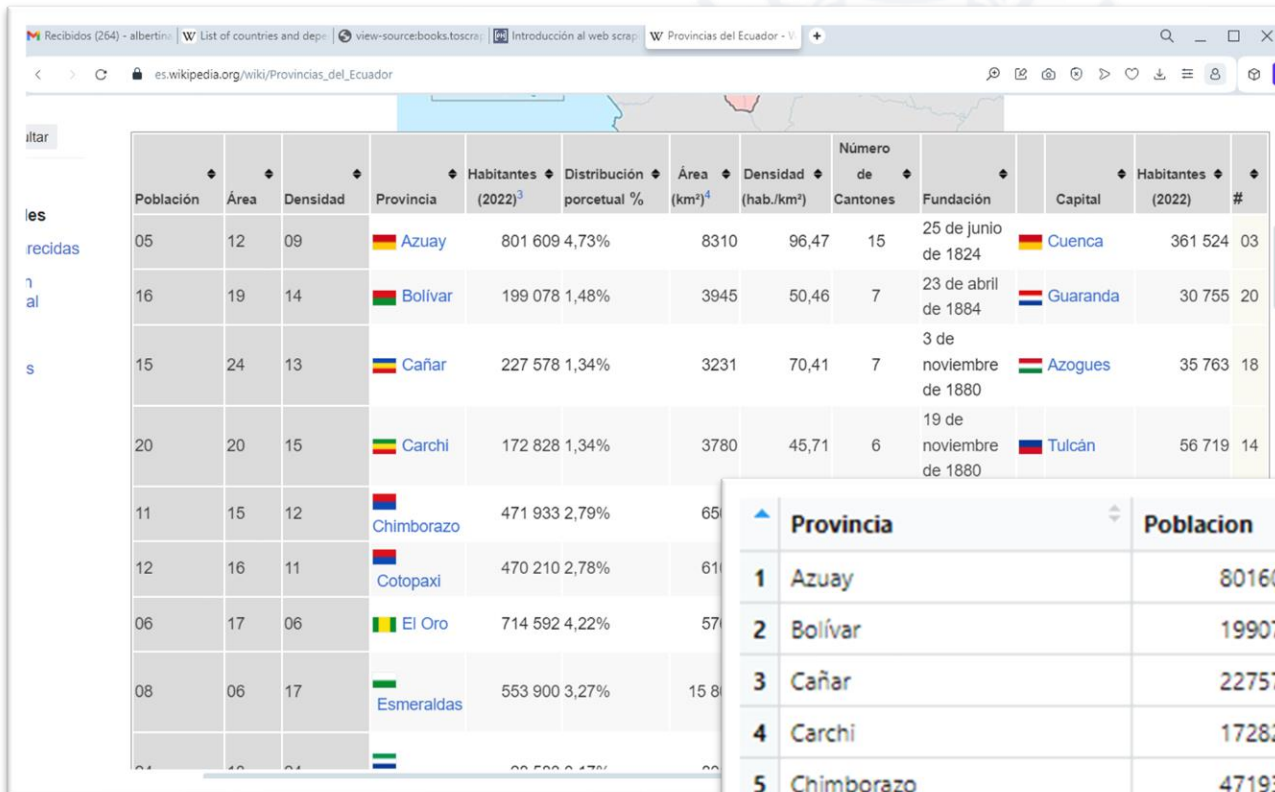
## Ejemplos de selectores CSS







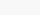
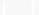
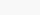



Selector CSS	Qué selecciona
h1	Todos los encabezados <h1>
.precio	Todos los elementos con clase precio
#producto123	El elemento con ID producto123
div > p	Todos los párrafos <p> que son hijos directos de un <div>
ul li a	Todos los enlaces dentro de listas
img[src\$=".jpg"]	Todas las imágenes cuyo atributo src termina en .jpg
a[href*="producto"]	Enlaces con href que contienen la palabra "producto"
div[class^="item"]	Todos los div cuya clase comienza con item

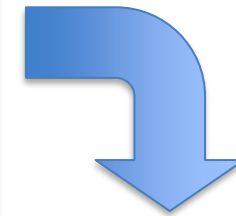
# Ejemplo Web Scrapping con R

[https://es.wikipedia.org/wiki/Provincias\\_del\\_Ecuador](https://es.wikipedia.org/wiki/Provincias_del_Ecuador)

## EJEMPLO 1



Población	Área	Densidad	Provincia	Habitantes (2022) <sup>3</sup>	Distribución porcentual %	Área (km²) <sup>4</sup>	Densidad (hab./km²)	Número de Cantones	Fundación	Capital	Habitantes (2022)	#
05	12	09	 Azuay	801 609	4,73%	8310	96,47	15	25 de junio de 1824	 Cuenca	361 524	03
16	19	14	 Bolívar	199 078	1,48%	3945	50,46	7	23 de abril de 1884	 Guaranda	30 755	20
15	24	13	 Cañar	227 578	1,34%	3231	70,41	7	3 de noviembre de 1880	 Azogues	35 763	18
20	20	15	 Carchi	172 828	1,34%	3780	45,71	6	19 de noviembre de 1880	 Tulcán	56 719	14
11	15	12	 Chimborazo	471 933	2,79%	6500	72,61	10				
12	16	11	 Cotopaxi	470 210	2,78%	6108	76,98	7				
06	17	06	 El Oro	714 592	4,22%	5767	123,92	14				
08	06	17	 Esmeraldas	553 900	3,27%	15809	35,04	7				



*Data frame PROVINCIAS*

	Provincia	Poblacion	Area	Densidad	Cantones
1	Azuay	801609	8310	96.47	15
2	Bolívar	199078	3945	50.46	7
3	Cañar	227578	3231	70.41	7
4	Carchi	172828	3780	45.71	6
5	Chimborazo	471933	6500	72.61	10
6	Cotopaxi	470210	6108	76.98	7
7	El Oro	714592	5767	123.92	14
8	Esmeraldas	553900	15809	35.04	7
9	Galápagos	28562	8010	3.57	2

# Extracción mediante web scraping

## EJEMPLO 1

```
1 # EJEMPLO WEB SCRAPING CON R
2 # (Usando librería rvest)
3
4 # Obtener información de las provincias del ecuador a partir de
5 # información disponible en Wikipedia
6
7 # Cargar las bibliotecas necesarias
8 install.packages("rvest")
9 install.packages("xml2")
10
11 library(rvest)
12 library(xml2)
13 library(dplyr)
14
15 # Definir la URL
16 url <- "https://es.wikipedia.org/wiki/Provincias_de_Ecuador"
17
18 # Leer el contenido de la página
19 pagina = read_html(url)
20
21 # Extraer tablas de la página
22 tablas = html_table(pagina)
23
24 # Asumiendo que la tabla de interés es la primera
25 # que contiene las Provincias de Ecuador
26 tabla_provincias = tablas[[1]]
```

# Extracción mediante web scraping

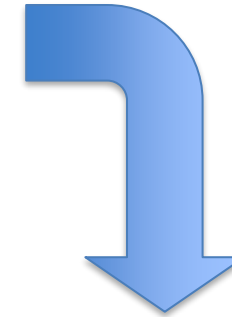
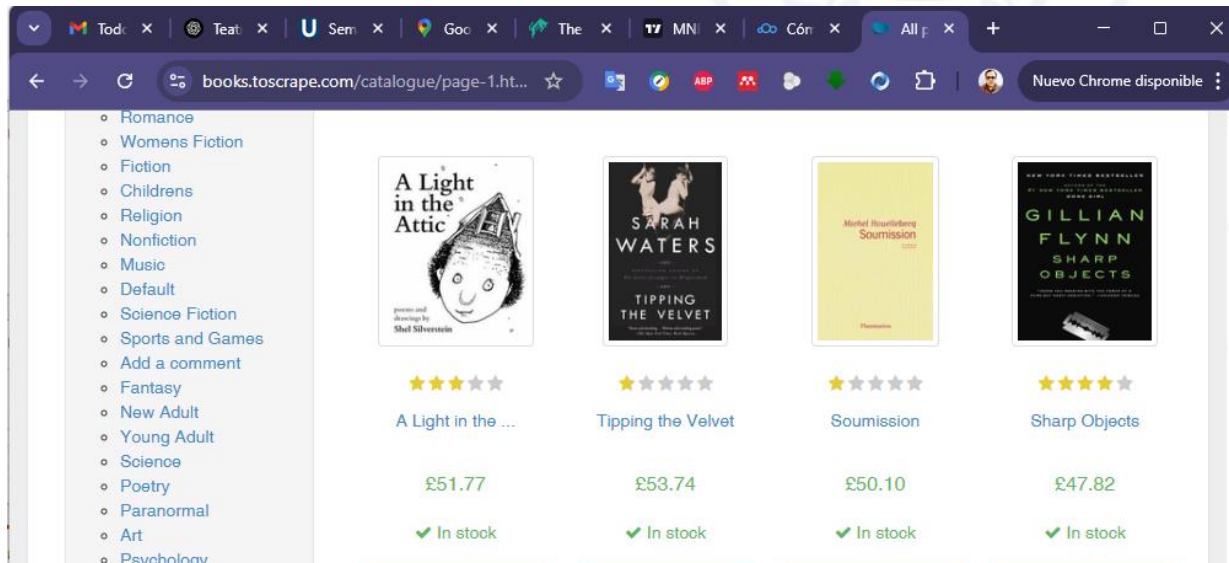
## EJEMPLO 1

```
28 # Inspeccionar la estructura del data frame
29 names(tabla_provincias)
30 str(tabla_provincias)
31
32 # Extraer únicamente información de nombre de la provincia,
33 # población, superficie, densidad poblacional, y número de cantones
34 provincias = tabla_provincias %>%
35   select(Provincia, Poblacion = 5, Area = 7, Densidad = 8, Cantones = 9)
36
37 str(provincias)
38
39 # Eliminar caracteres distintos a dígitos en columnas numéricas enteras,
40 # y reemplazar la coma decimal por punto en columnas numéricas decimales,
41 # y convertir a valores enteros o decimales.
42 provincias = provincias %>%
43   mutate(Poblacion = as.integer(gsub("[^0-9]", "", Poblacion)),
44          Area = as.integer(gsub("[^0-9]", "", Area)),
45          Densidad = as.numeric(gsub(",", ".", Densidad)))
46
47 str(provincias)
48 print(provincias)
```

# Ejemplo Web Scrapping con R

<https://books.toscrape.com/catalogue/page-1.html>

## EJEMPLO 2



*Data frame LIBROS*

	titulo	enlace	precio
1	A Light in the Attic	a-light-in-the-attic_1000/index.html	51.77
2	Tipping the Velvet	tipping-the-velvet_999/index.html	53.74
3	Soumission	soumission_998/index.html	50.10
4	Sharp Objects	sharp-objects_997/index.html	47.82
5	Sapiens: A Brief History of Humankind	sapiens-a-brief-history-of-humankind_996/index.html	54.23
6	The Requiem Red	the-requiem-red_995/index.html	22.65
7	The Dirty Little Secrets of Getting Your Dream Job	the-dirty-little-secrets-of-getting-your-dream-job_994/inde...	33.34
8	The Coming Woman: A Novel Based on the Life of the Infam...	the-coming-woman-a-novel-based-on-the-life-of-the-infa...	17.93

# Extracción mediante web scraping

## EJEMPLO 2

```
1 # Obtener lista de títulos y enlaces de libros publicados
2 # en el sitio books.toscrape.com
3
4 # Cargar librerías necesarias
5 library(rvest)
6 library(xml2)
7
8 # Leer el contenido de la página web
9 url <- "https://books.toscrape.com/catalogue/page-1.html"
10 pagina <- read_html(url)
11
12 # Extraer información de los contenedores que contienen la información de los
13 # libros utilizando selectores CSS
14 nodos = pagina %>%
15   html_elements(".product_pod")
16
17 # Extraer títulos, enlaces, precio, e imágenes de los libros
18 titulos = nodos %>%
19   html_element("h3 a") %>%
20   html_attr("title")
21
22 enlaces = nodos %>%
23   html_element("h3 a") %>%
24   html_attr("href")
```

# Extracción mediante web scraping

## EJEMPLO 2

```
26 precios = nodos %>%
27   html_element(".price_color") %>%
28   html_text2()
29
30 imagenes = nodos %>%
31   html_element(".image_container a img") %>%
32   html_attr("src")
33
34 imagenes = gsub("\\.\\.\\.\"", "", imagenes)
35
36 # Combinar vectores en un solo data frame
37 libros = data.frame(titulo = titulos,
38                     enlace = paste0("http://books.toscrape.com/catalogue/", enlaces),
39                     precio = precios,
40                     imagen = paste0("http://books.toscrape.com", imagenes))
41
42 # Corregir el formato del precio de los libros
43 libros = libros %>%
44   mutate(precio = as.numeric(gsub("£", "", precio)))
45 str(libros)
46
47 # Inspeccionar el data frame final resultante
48 print(libros)
49
```

UNIVERSIDAD TÉCNICA PARTICULAR DE LOJA

# Extracción mediante web scraping

- Recursos recomendados
  - Web scraping in R (<https://www.r-bloggers.com/2023/01/web-scraping-in-r-2/>)
- Revisar más acerca de:
  - CSS selectors
  - XPath



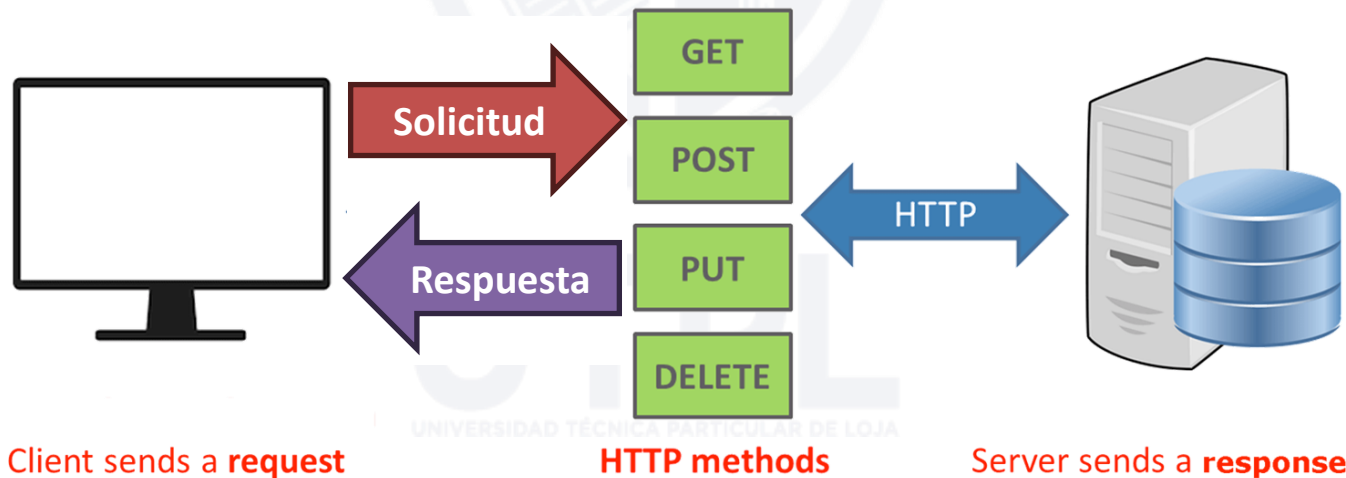


# Ejercicio Web Scraping

- Ir al sitio de Alpha Editorial (<https://www.alphaeditorialcloud.com/>)
- Realizar una búsqueda de libros bajo cualquier criterio.
- Del resultado de la búsqueda extraer:
  - Título del libro
  - Fecha
  - Autor
  - Precio
  - URL de la imagen portada del libro

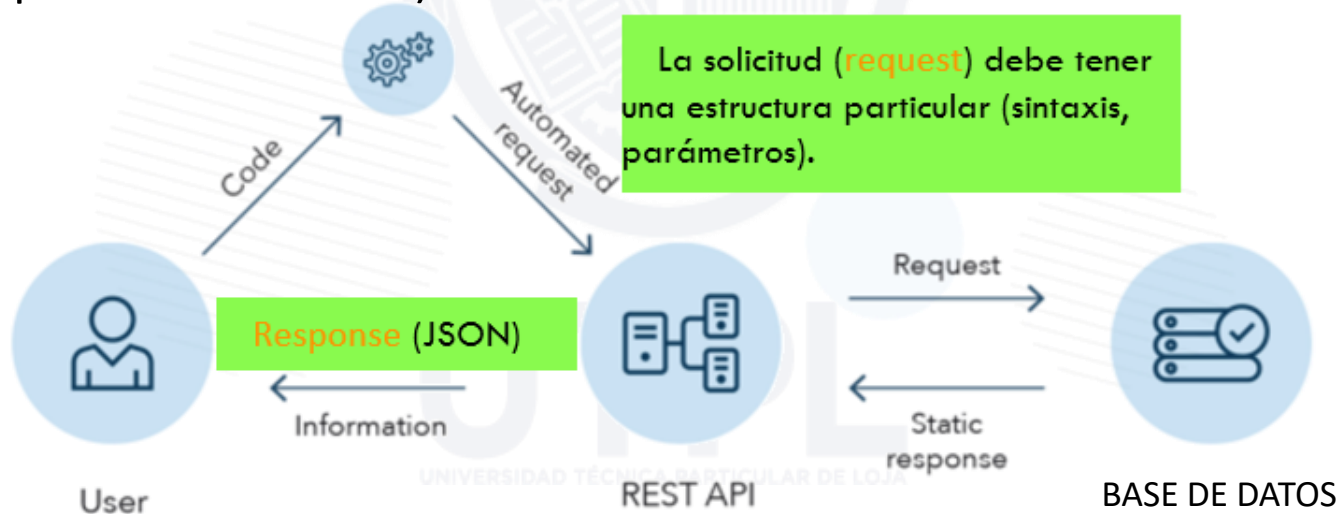
# Extracción mediante APIs

- **API** (Application Programming Interface) es la interfaz que permite la comunicación (intercambio de datos) entre dos aplicaciones de software independientes a través de un conjunto de reglas,
- Una API se puede implementar mediante un **servicio web** que otorga acceso a datos y métodos específicos a los que otras aplicaciones pueden acceder, y en ocasiones editar, a través de protocolos **HTTP** estándar.



# Extracción mediante APIs

- **REST** (REpresentational State Transfer) es el estilo arquitectónico más popular de API para servicios web que ofrecen los proveedores de datos.
- **REST** se refiere a un conjunto de pautas diseñadas para simplificar la comunicación cliente/servidor. Las API REST hacen que el acceso a los datos sea mucho más sencillo y lógico.
- Pueden ser privadas (empresa), abiertas sólo para partners, o públicas (cualquier desarrollador).



# Extracción mediante APIs

¿Qué es una API REST o API RESTful?



<https://youtu.be/gUeyJ9xDq-Y?si=PcZ9WIfFsPF9APJu>

REST => R (read)

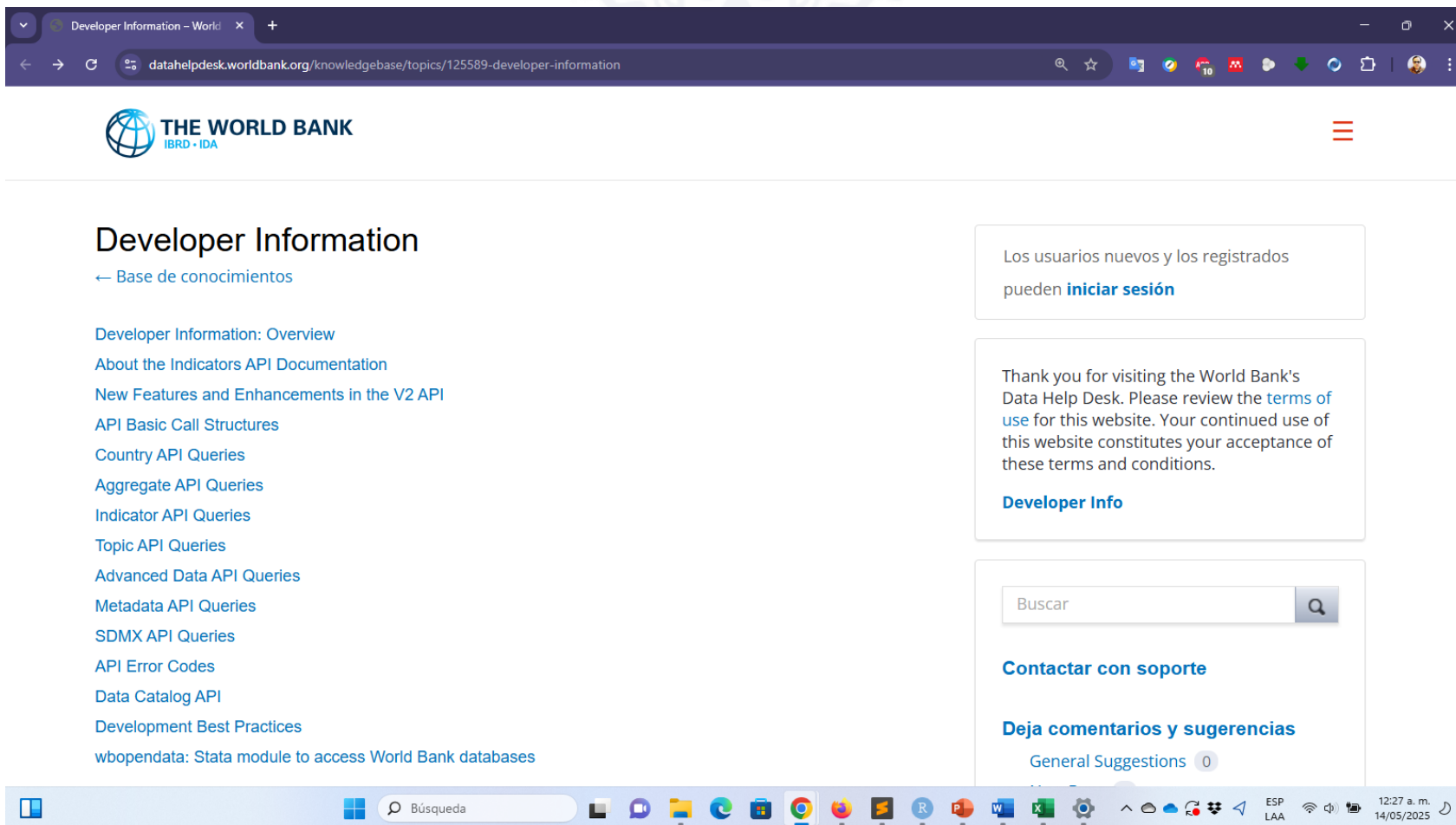
RESTful => CRUD (Create, Read, Update, Delete)

# Extracción mediante APIs

- En R se utiliza la librería “**httr**” que facilita la comunicación con servicios API RESTful.
- También se requiere usar la librería “**jsonlite**” que permite trabajar con archivos JSON, que es el formato comúnmente usado para recepción o envío de datos a través de APIs.
- Se requiere revisar la documentación de la API para conocer su tipo, restricciones, parámetros, estructura del REQUEST, etc.
- Ejemplo:
  - API: Banco Mundial
  - EndPoint: <http://api.worldbank.org/V2>
  - Documentación: <https://datahelpdesk.worldbank.org/knowledgebase/topics/125589>

# Extracción mediante APIs

## Documentación



The screenshot shows a web browser window displaying the 'Developer Information' page of the World Bank Data Help Desk. The browser's address bar shows the URL: `datahelpdesk.worldbank.org/knowledgebase/topics/125589-developer-information`. The page features the World Bank logo (IBRD - IDA) and a navigation menu on the left with links such as 'Developer Information: Overview', 'About the Indicators API Documentation', 'New Features and Enhancements in the V2 API', 'API Basic Call Structures', 'Country API Queries', 'Aggregate API Queries', 'Indicator API Queries', 'Topic API Queries', 'Advanced Data API Queries', 'Metadata API Queries', 'SDMX API Queries', 'API Error Codes', 'Data Catalog API', 'Development Best Practices', and 'wbopendata: Stata module to access World Bank databases'. On the right side, there are three informational boxes: the first mentions login for new and registered users; the second is a thank-you message with a link to the 'terms of use'; the third contains a search bar, a 'Contactar con soporte' link, and a 'Deja comentarios y sugerencias' link with a 'General Suggestions' counter showing 0. The Windows taskbar at the bottom includes a search bar, various application icons, and system status information indicating the time as 12:27 a.m. on 14/05/2025.

Developer Information  
← Base de conocimientos

Developer Information: Overview  
About the Indicators API Documentation  
New Features and Enhancements in the V2 API  
API Basic Call Structures  
Country API Queries  
Aggregate API Queries  
Indicator API Queries  
Topic API Queries  
Advanced Data API Queries  
Metadata API Queries  
SDMX API Queries  
API Error Codes  
Data Catalog API  
Development Best Practices  
wbopendata: Stata module to access World Bank databases

Los usuarios nuevos y los registrados pueden [iniciar sesión](#)

Thank you for visiting the World Bank's Data Help Desk. Please review the [terms of use](#) for this website. Your continued use of this website constitutes your acceptance of these terms and conditions.

[Developer Info](#)

Buscar

[Contactar con soporte](#)

[Deja comentarios y sugerencias](#)  
General Suggestions 0

# Extracción mediante APIs

## Documentación

The screenshot shows a web browser window displaying the 'API Basic Call Structures' page from the World Bank's Data Help Desk. The page is titled 'API Basic Call Structures' with a sub-link 'Developer Information'. The main content explains that the Indicators API supports two basic ways to build queries: a URL based structure and an argument based structure. It provides two examples of URLs to retrieve data for low-income countries. Below this, the 'Query Strings' section explains that the API supports query strings in requests, with examples for date ranges and ranges. The right sidebar contains a search bar, a 'Contactar con soporte' link, a 'Deja comentarios y sugerencias' section with counts for General Suggestions, New Data, and Website Improvements, and a 'Base de conocimientos' section with counts for Country Classification, Currencies, Data Compilation Methodology, Data Not Available, Data Updates, and DataBank. The bottom of the browser window shows the Windows taskbar with various application icons and the system clock indicating 12:29 a.m. on 14/05/2025.

API Basic Call Structures  
← Developer Information

The Indicators API supports two basic ways to build queries: a URL based structure and an argument based structure. For example, the following two requests will return the same data, a list of countries with income level classified as low income:

- Argument based: <https://api.worldbank.org/V2/country?incomeLevel=LIC>
- URL based: <https://api.worldbank.org/V2/incomeLevel/LIC/country>

### Query Strings

The API support the following query strings in requests.

**Date and Date-Range:** Date-range by year, month or quarter that scopes the result-set.

Examples:

- <https://api.worldbank.org/v2/country/all/indicator/SP.POP.TOTL?date=2000>
- <https://api.worldbank.org/v2/country/chn/bra/indicator/DPANUSSPB?date=2012M01>

A range is indicated using the colon ( : ) separator.

Examples:

- <https://api.worldbank.org/v2/country/all/indicator/SP.POP.TOTL?date=2000:2001>

Los usuarios nuevos y los registrados pueden [iniciar sesión](#)

Thank you for visiting the World Bank's Data Help Desk. Please review the [terms of use](#) for this website. Your continued use of this website constitutes your acceptance of these terms and conditions.

[Developer Info](#)

Buscar

[Contactar con soporte](#)

[Deja comentarios y sugerencias](#)

General Suggestions 0  
New Data 0  
Website Improvements 0

[Base de conocimientos](#)

Country Classification 4  
Currencies 3  
Data Compilation Methodology 12  
Data Not Available 7  
Data Updates 3  
DataBank 13

# Extracción mediante APIs

- Pasos
  1. Construir la URL de la llamada a la API
  2. Realizar una solicitud GET()
  3. Verificar el estado de la respuesta:
    - Se espera siempre un status 200-OK
    - Otros status pueden ser: 400-Bad Request, 401-Unauthorized, 403-Forbidden, 404-Not Found
  4. Traducir la respuesta a un formato utilizable, como un data frame.



# Extracción mediante APIs

- **GET()** genera un objeto ***response***, que contiene toda la información del *request*, que incluye:
  - La **url** a la que finalmente se envió la solicitud
  - El **status\_code** HTTP
  - Una referencia interna (**handle**) usada por *httr* para gestionar la solicitud.
  - Los encabezados HTTP (**headers**) de la respuesta. Por ejemplo, tipo de contenido, permisos, etc.
  - El contenido o cuerpo de la respuesta (**content**), en un formato determinado (XML, JSON, etc.)
  - Información fecha y hora (**time**) de la solicitud
  - ...

# Extracción mediante APIs

## EJEMPLO

- Usando la API del Banco Mundial (<http://api.worldbank.org/>), # extraer información países del mundo
- URL base: <http://api.worldbank.org/v2/country>
- Parámetros principales del REQUEST: **region** (región geográfica: LCN, EAS, ECS, ...), **incomeLevel** (nivel de ingresos: LIC, LMC, UMC, HIC), **format** (json, xml), **per\_page** (cantidad observaciones).
- Documentación:
  - <https://datahelpdesk.worldbank.org/knowledgebase/articles/898590-country-api-queries>
  - <https://api.worldbank.org/v2/region?format=json>
  - <https://api.worldbank.org/v2/incomeLevel?format=json>

# Extracción mediante APIs

## EJEMPLO

```
4 # Usando la API del Banco Mundial (http://api.worldbank.org/),
5 # extraer información de todos los países del mundo
6
7 # Instalar y cargar librerías
8 install.packages("httr")
9
10 library(httr)
11 library(jsonlite)
12 library(dplyr)
13
14 # Definir la URL base y parámetros para la API del World Bank
15 url_base <- "http://api.worldbank.org/v2/country"
16
17 region <- "CLA"
18 nivel_ing <- "UMC"
19 formato <- "json"
20 cant_obs <- 1000
21
22 # Construir la URL completa
23 url <- paste0(url_base,
24               "?region=", region,
25               "&incomeLevel=", nivel_ing,
26               "&format=", formato,
27               "&per_page=", cant_obs)
```

URL: [http://api.worldbank.org/v2/country?region=CLA&incomeLevel=UMC&format=json&per\\_page=100](http://api.worldbank.org/v2/country?region=CLA&incomeLevel=UMC&format=json&per_page=100)

# Extracción mediante APIs

## EJEMPLO

```
29 # Realizar la solicitud GET
30 respuesta <- GET(url)
31 class(respuesta)
32
33 # Verificar el estado de la respuesta (se espera status HTTP 200)
34 status_code(respuesta)
35
36 # Si la solicitud fue exitosa, extraer y procesar los datos
37
38 # Extraer el contenido del response (en formato JSON)
39 contenido_json <- content(respuesta,"text")
40
41 contenido_R <- fromJSON(contenido_json, flatten = TRUE)
42
43 # Se extrae el 2do elemento de la lista que contiene el data frame
44 # con los datos solicitados.
45 paises <- contenido_R[[2]]
46
47 # Mostrar los datos
48 print(paises)
```

# Extracción usando APIs

## EJEMPLO

```
[
  {
    page: 1,
    pages: 1,
    per_page: "1000",
    total: 18
  },
  [
    {
      id: "ARG",
      iso2Code: "AR",
      name: "Argentina",
      region: {
        id: "LCN",
        iso2code: "ZJ",
        value: "Latin America & Carib
      },
      adminregion: {
        id: "LAC",
        iso2code: "XJ",
        value: "Latin America & Carib
      },
      incomeLevel: {
        id: "UMC",
        iso2code: "XT",
        value: "Upper middle income"
      },
      lendingType: {
        id: "IBD",
        iso2code: "XF",
        value: "IBRD"
      },
      capitalCity: "Buenos Aires",
      longitude: "-58.4173",
      latitude: "-34.6118"
    }
  ]
]
```



	id	iso2Code	name	capitalCity	longitude
1	ARG	AR	Argentina	Buenos Aires	-58.4173
2	BLZ	BZ	Belize	Belmopan	-88.7713
3	BRA	BR	Brazil	Brasilia	-47.9292
4	COL	CO	Colombia	Bogota	-74.082
5	CRI	CR	Costa Rica	San Jose	-84.0089
6	DMA	DM	Dominica	Roseau	-61.39
7	DOM	DO	Dominican Republic	Santo Domingo	-69.8908
8	ECU	EC	Ecuador	Quito	-78.5243
9	GRD	GD	Grenada	Saint George's	-61.7449
10	GTM	GT	Guatemala	Guatemala City	-90.5328
11	JAM	JM	Jamaica	Kingston	-76.792
12	LCA	LC	St. Lucia	Castries	-60.9832
13	MEX	MX	Mexico	Mexico City	-99.1276
14	PER	PE	Peru	Lima	-77.0465
15	PRY	PY	Paraguay	Asuncion	-57.6362
16	SLV	SV	El Salvador	San Salvador	-89.2073

# Extracción usando APIs

- Lista de APIs públicas:

<https://github.com/public-apis/public-apis>

# Ejercicio extracción con APIs

- Usando la API que proporciona la plataforma OpenWeather (<https://api.openweathermap.org/>): Extraer datos del pronóstico del tiempo de los próximos 5 días, para un lugar específico.
- URL base: <https://api.openweathermap.org/data/2.5/forecast>
- Documentación de la API: <https://openweathermap.org/forecast5>
- La API requiere una llave (API-Key), usar:  
**4afd29be646ee3c7027bea84dc25ce84**
- Usar el menos 4 parámetros y extraer información correspondiente a la hora, temperatura mínima, temperatura máxima, nivel de lluvia, y tipo de clima.