



UT03: GESTIÓN DE LOS DATOS

ÍNDICE

- 1.- Introducción a ETL con Pandas
- 2.- Extracción: conectando con las fuentes
- 3.- Transformación: el corazón del proceso
- 4.- Carga: almacenamiento y salida
- 5.- Optimización y buenas prácticas
- 6.- Proceso ETL en entorno gráfico: Apache NiFi

1

INTRODUCCIÓN A LOS PIPELINES DE DATOS



1

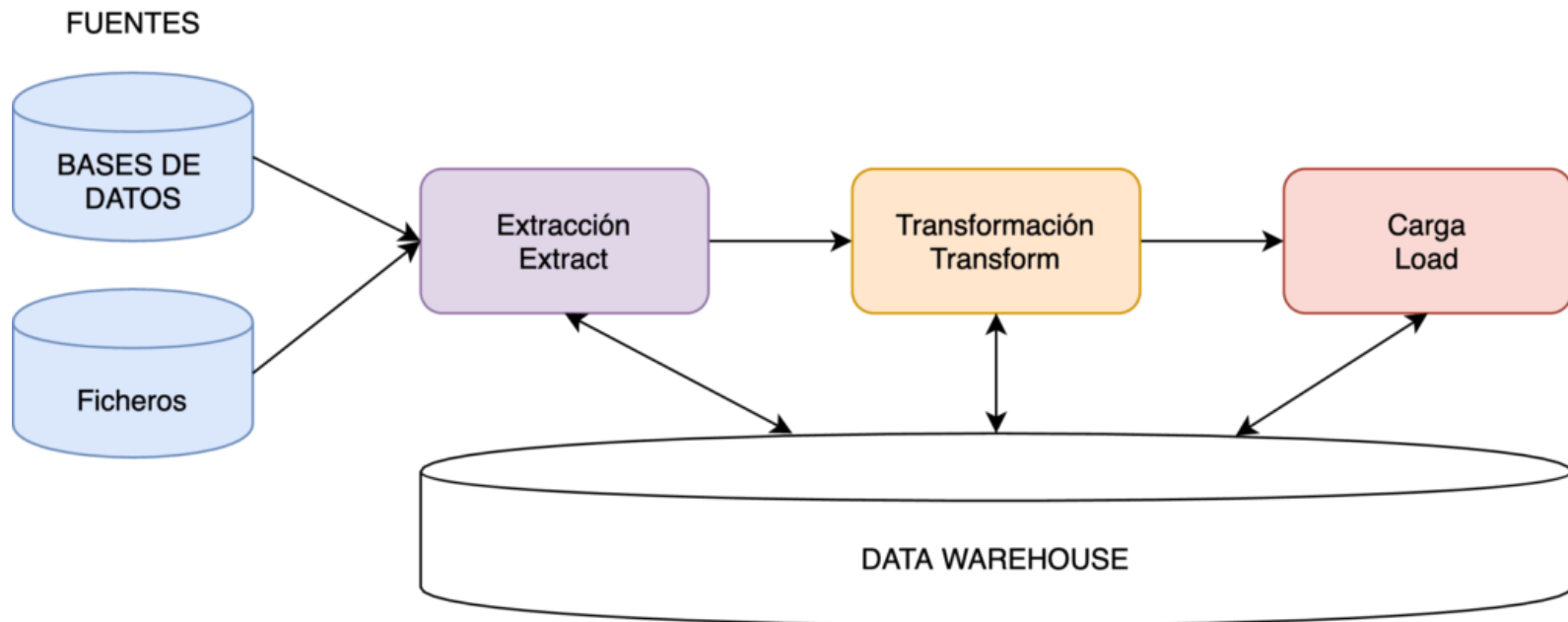
INTRODUCCIÓN A LOS PIPELINES DE DATOS

1.1

¿QUÉ ES UN PIPELINE DE DATOS?

Un **pipeline de datos** es un conjunto de procesos estructurados que permiten **mover, transformar y almacenar datos** desde su origen hasta su destino final.

Este flujo suele automatizarse y definirse de forma secuencial.

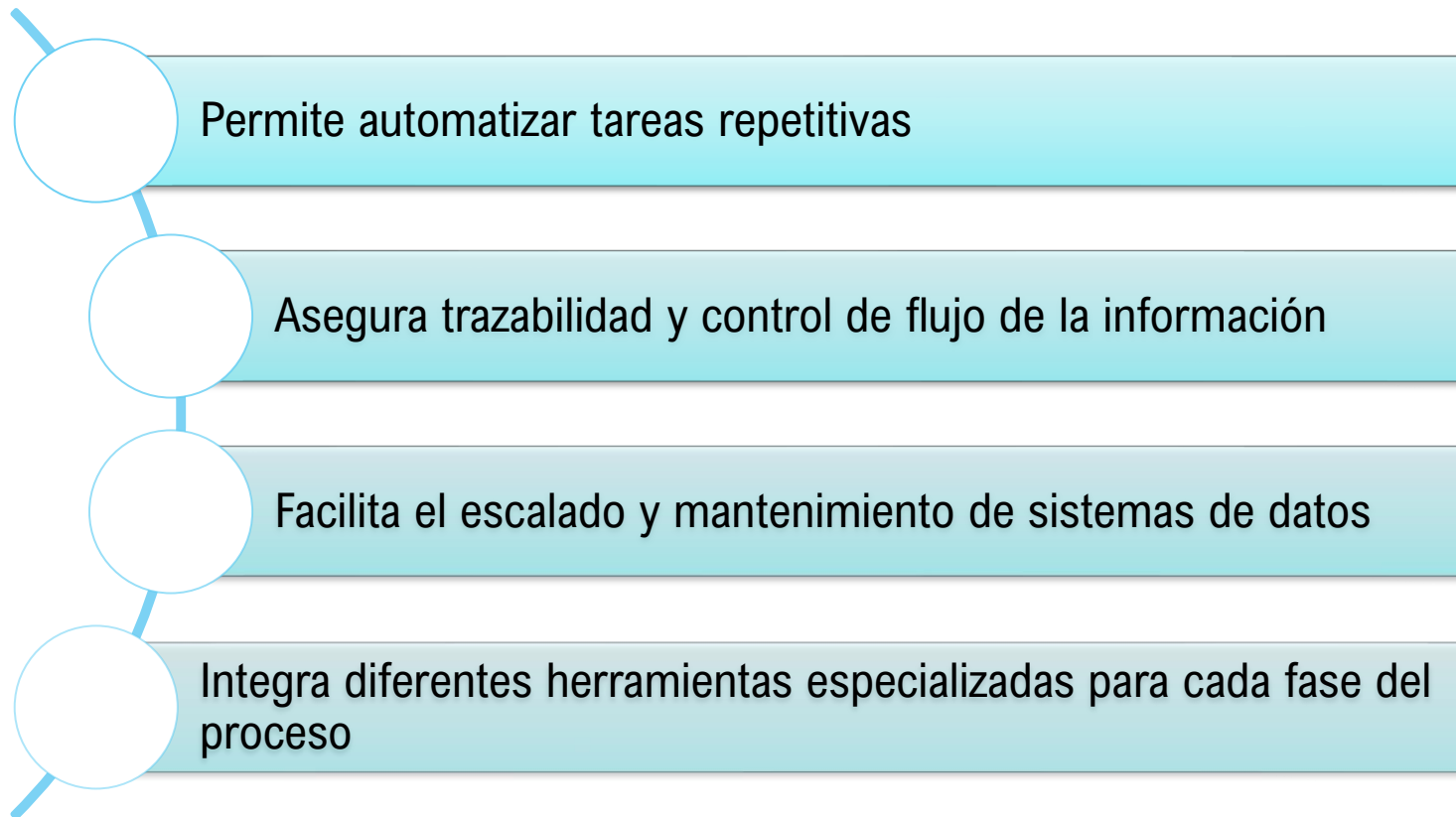


El concepto no se refiere solo a mover los datos, sino también a **transformarlos y prepararlos** a lo largo del proceso, garantizando su calidad, integridad y utilidad.

Algunas etapas por las que pasa son:

- Limpieza
- Validación
- Enriquecimiento
- Combinación con otras fuentes

Un pipeline de datos bien diseñado:



Un pipeline de datos se compone de varias **fases encadenadas**. Estas fases son:

- Ingesta de datos
- Procesamiento ETL y ELT
- Almacenamiento de datos
- Análisis de datos
- Consumo de datos

1. Ingesta de datos

La ingesta consiste en **capturar y recoger datos desde múltiples fuentes**, que pueden ser estructuradas (BBDD SQL), semiestructuradas (JSON, XML) o no estructuradas (log, imágenes, texto)

Ejemplos de fuentes:

- APIs (REST, SOAP)
- Bases de datos SQL y NoSQL
- Archivos planos (CSV, Excel, Parquet)
- Logs de servidores
- Sensores IoT
- Servicios de streaming (Kafka, MQTT)

La ingesta puede ser por lotes periódicos (**batch**) o en tiempo real (**streaming**)

2. Procesamiento: ETL y ELT

Esta fase se encarga de **preparar la información para su análisis y consumo.**

Hay dos enfoques diferentes:

- **ETL:** Extract, Transform, Load
- **ELT:** Extract, Load, Transform

ETL

Extract – Transform – Load

Es el enfoque utilizado en arquitecturas de **Data Warehousing**

Los pasos son:

- **Extract:** se recopilan datos desde diferentes fuentes
- **Transform:** se **aplica lógica de negocio** a los datos. Por ejemplo:
 - Tareas de limpieza (eliminación de duplicados, nulos)
 - Conversión de tipos de datos
 - Normalización y estandarizaciones
 - Agregaciones y cálculos
 - Enriquecimiento
- **Load:** finalmente los datos se almacenan en un sistema destino **estructurado** como un Data Warehouse

ELT

Extract – Load - Transform

Es el enfoque utilizado en arquitecturas de **Cloud Data Warehousing** y **Data Lakes**

Los pasos son:

- **Extract:** se recopilan datos desde diferentes fuentes
- **Load:** los datos se cargan inmediatamente en el sistema destino en su estado original (*raw*) sin modificaciones previas
- **Transform:** se aplica la lógica a los datos dentro del sistema destino, por ejemplo:
 - Ejecución de consultas SQL para limpiar y filtrar datos cargados
 - Modelado de datos y creación de vistas para análisis
 - Generación de tablas agregadas para reportes finales

1

INTRODUCCIÓN A LOS PIPELINES DE DATOS

1.2

INTRODUCCIÓN A PANDAS

¿Qué es Pandas?

Pandas es una biblioteca de software de código abierto escrita para el lenguaje de programación Python.

Su propósito fundamental es la **manipulación y el análisis de datos**.

Está construido sobre **Numpy**, otra librería de Python especializada en cálculos numéricos sobre arrays de datos.



Pandas, cuyo nombre viene PANel DAta, fue desarrollado en 2008 por **Wes McKinney** mientras trabajaba en AQR Capital Management.

En 2009 fue liberado como Open Source.

En la actualidad es el **estándar de facto** en la manipulación de datos programática.

2

EXTRACCIÓN: CONECTANDO CON LAS FUENTES



2

EXTRACCIÓN: CONECTANDO CON LAS FUENTES

2.1

FORMATOS DE ARCHIVOS COMUNES: CSV

Algunos ejemplos de archivos comunes que podemos encontrar son:

- Archivos de texto plano
- Ficheros CSV (*Comma Separated Values*)
- Documentos Excel
- Archivos JSON
- Formatos específicos: Parquet y Feather

Archivos CSV

Los archivos **CSV (Comma Separated Values)** son un formato universal para el intercambio de datos.

Son simples, legibles por humanos y soportados por casi cualquier sistema informático, desde mainframes de los años 80 hasta las nubes modernas.

Aunque comúnmente los llamamos CSV, en realidad son **Archivos de Texto Delimitados**, y tienen 3 componentes clave que debemos tratar:

- El delimitador
- La codificación
- El encabezado

El delimitador

La codificación

El encabezado

Es el carácter que separa una columna de la otra.

Puede ser:

- **Coma (,)**: el estándar oficial del CSV
- **Punto y coma (;)**: muy común en Europa y Latinoamérica, ya que son lugares donde se usa la coma como separador decimal (10,5)
- **Tabulador(\t)**: se llaman archivos TSV (Tab Separated Values). Son muy seguros porque es raro que un archivo contenga una tabulación accidentalmente.
- **Pipe (|)**: usado a menudo en sistemas antiguos o Unix
- **Delimitadores multicarácter**: a veces, un único carácter no es seguro porque el texto mismo podría contenerlo, por lo que en ocasiones podemos encontrar secuencias de dos o más caracteres como :: o ~!~

- **Espacios en blanco variables:** es el clásico en archivos de logs de servidores o salidas de comando de Linux, donde hay uno o más espacios de forma que se vean visualmente alineados.

IP	USER	STATUS
192.168.1.1	admin	OK
10.0.0.5	guest	ERROR

- **Caracteres invisibles (ASCII Control Characters):** en el mundo de Hadoop, Hive y Mainframes se suelen usar caracteres que no se pueden imprimir ni ver en un editor normal. Esto garantiza al 100% que el delimitador nunca aparezca dentro del texto de un usuario. Por ejemplo, `\x01` Start of Heading (es el delimitador por defecto de Hive) o `\0` (Null byte)

El delimitador

La codificación

El encabezado

Es la tabla de caracteres que usa el archivo para entender letras y símbolos. La elección de un sistema de codificación incorrecto mostrará caracteres extraños (Españ@a en lugar de España).

Los sistemas más habituales son:

- **UTF-8:** el estándar moderno. Soporta emojis y todos los idiomas. Siempre intenta este primero.
- **Latin-1 (ISO-8859-1):** el estándar antiguo de Windows/Europa Occidental. Muy común en archivos generados por Excel antiguos o sistemas *legacy* bancarios

El delimitador

La codificación

El encabezado

Opcionalmente, la primera fila del archivo puede contener los nombres de las columnas.

Si el fichero no tiene encabezado, habrá que proporcionar los nombres manualmente durante la extracción.

pd.read_csv()

Carga un fichero CSV en un dataframe

```
import pandas as pd
```

```
df = pd.read_csv('datos.csv')
```

```
print(df.head())
```

```
print(df.dtypes)
```

Pandas detecta automáticamente encabezados y separadores

	ID	Nombre	Departamento	Salario_Anual	Fecha_Contratacion
0	1001	Ana García	Marketing	45000	2021-03-15
1	1002	Carlos Ruiz	Ventas	38500	2022-07-01
2	1003	Elena Vázquez	Desarrollo	52000	2020-11-20
3	1004	Jorge Méndez	Recursos Humanos	32000	2023-01-10
4	1005	Lucía Torres	Finanzas	47000	2019-05-25

ID	int64
Nombre	object
Departamento	object
Salario_Anual	int64
Fecha_Contratacion	object
dtype:	object

También reconoce los tipos de datos de las columnas

Algunos parámetros que puede tener son:

- **sep** o **delimiter**: el carácter que separa las columnas, por defecto coma
- **header**: fila que contiene los nombres de las columnas. 0 es la primera y None para indicar que no tiene.
- **names**: si header=None o se quieren renombrar. Es una lista con una cadena por cada campo.
- **index_col**: usa una columna del archivo como índice etiqueta de fila en lugar de números consecutivos. Ej.: `index_col='ID_Cliente'`
- **use_cols**: carga solo las columnas que se indiquen. Ahorra mucha RAM
- **nrows**: lee solo las N primeras filas. Útil para pruebas
- **skiprows**: salta filas al inicio (títulos, ...)

- **dtype:** fuerza el tipo de dato de las columnas. Ej.: `dtype={'ID': str, 'Edad': 'int32'}`
- **parse_dates:** intenta convertir columnas a objetos `datetime` automáticamente. Ej.: `parse_dates=['fecha_compra']`
- **converters:** aplica una función de Python a los datos mientras lee. Lento pero potente. Ej.: `converters={'precio': limpiar_moneda}`
- **na_values:** qué hacer si una fila tiene más columnas de las esperadas. Los posibles valores son `error` (opción por defecto), `skip` y `warn`
- **encoding:** codificación de caracteres, algunos valores pueden ser `utf-8`, `latín-1` o `cp1252`

Ejemplo:

```
import pandas as pd

df = pd.read_csv(
    'datos.csv',
    sep=',',
    header=0,
    index_col='ID',
    encoding='utf-8',
    parse_dates=['Fecha_Contratacion'],
    dtype={'Departamento': 'category'}
)

print(df.info())
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 5 entries, 1001 to 1005
Data columns (total 4 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Nombre                5 non-null     object
1   Departamento           5 non-null     category
2   Salario_Anual          5 non-null     int64
3   Fecha_Contratacion     5 non-null     datetime64[ns]
dtypes: category(1), datetime64[ns](1), int64(1), object(1)
memory usage: 377.0+ bytes
None
```

Aunque la función `read_csv()` funciona muy bien, **en el mundo real lo habitual es encontrar archivos mal formados o corruptos** que hacen que el script lance un error fatal a mitad de proceso.

Veamos cómo gestionar algunos de los **errores más comunes** al leer archivos de texto con Pandas.

Error de estructura (`ParseError`)

Ocurre cuando una fila tiene más o menos columnas de las que debería. Por defecto, Pandas lanzará un `ParseError` y detendrá todo.

La solución es utilizar el parámetro `on_bad_lines='warn'`, que mostrará un aviso pero no interrumpirá la carga

Error de tipo de dato

Ocurre cuando una columna numérica tiene basura textual.

La estrategia es no forzar `dtype='int'`. Es mejor leer como `object` y limpiar en el paso de transformación.

	ID	Nombre	Departamento	Salario_Anual	Fecha_Contratacion
1	1001	Ana García	Marketing	45000	2021-03-15
2	1002	Carlos Ruiz	Ventas	38500	2022-07-01
3	1003	Elena Vázquez	Desarrollo	52000	2020-11-20
4	1004	Jorge Méndez	Recursos Humanos	ERROR	2023-01-10
5	1005	Lucía Torres	Finanzas	47000	2019-05-25

La estrategia es no forzar `dtype='int'` ya que lanzará error.

```
df = pd.read_csv(  
    'datos.csv',  
    sep=',',  
    encoding='utf-8',  
    dtype={'Salario_Anual': int}  
)
```

Si forzamos leerlo
como entero
lanzará error

```
print(df.head())
```

```
-----  
TypeError                                Traceback (most recent call last)  
File parsers.pyx:1160, in pandas._libs.parsers.TextReader._convert_tokens()
```

Es mejor leer como `object` y limpiar en el paso de transformación.

Leemos todo como cadenas

Convertimos a valor numérico

```
df = pd.read_csv('datos.csv', dtype=str)

df['Salario_Anual'] = pd.to_numeric(df['Salario_Anual'], errors='coerce')

datos_limpios.head()
```

Las posibilidades son:

raise: si algo no es número lanza error

ignore: si algo no es número deja la columna como cadena.

coerce: fuerza la conversión, si algo no es número convierte a NaN

	ID	Nombre	Departamento	Salario_Anual	Fecha_Contratacion
0	1001	Ana García	Marketing	45000.0	2021-03-15
1	1002	Carlos Ruiz	Ventas	38500.0	2022-07-01
2	1003	Elena Vázquez	Desarrollo	52000.0	2020-11-20
3	1004	Jorge Méndez	Recursos Humanos	NaN	2023-01-10
4	1005	Lucía Torres	Finanzas	47000.0	2019-05-25

Por ejemplo, eliminando las columnas con NaN

```
errores = df[df['Salario_Anual'].isna()]
datos_limpios = df.dropna(subset=['Salario_Anual'])

datos_limpios.head()
```

	ID	Nombre	Departamento	Salario_Anual	Fecha_Contratacion
0	1001	Ana García	Marketing	45000.0	2021-03-15
1	1002	Carlos Ruiz	Ventas	38500.0	2022-07-01
2	1003	Elena Vázquez	Desarrollo	52000.0	2020-11-20
4	1005	Lucía Torres	Finanzas	47000.0	2019-05-25

Valores nulos no estándar

Por defecto, Pandas sabe que NaN, null o una celda vacía son nulos.

Pero podemos encontrar con datos de sistemas antiguos a veces usan cosas creativas como ?, -, SIN_DATO o 9999.

La solución es utilizar el parámetro `na_values` al que se le pasa una lista que define a qué se le asigna el valor NaN en el archivo.

```
df = pd.read_csv(  
    'daots.csv',  
    na_values=['?', 'SIN_RESPUESTA', '-', 'n/a']  
)
```

El motor de lectura (engine)

A veces, el delimitador es complejo (ej: regex) o el archivo tiene caracteres muy extraños y el motor estándar de Pandas (escrito en C para velocidad) falla.

La solución es cambiar al motor de Python. Es más lento, pero mucho **más tolerante** y potente para casos extremos.

```
df = pd.read_csv(  
    'archivo_raro.txt',  
    sep='--',  
    engine='python')
```

Por ejemplo, el motor C falla bastante con separadores multicarácter

2

EXTRACCIÓN: CONECTANDO CON LAS FUENTES

2.2

FORMATOS DE ARCHIVOS COMUNES: EXCEL

Leer archivos de Excel es diferente a leer CSV.

Mientras que CSV es texto plano simple, un archivo Excel (.xlsx) es técnicamente un contenedor comprimido (ZIP) lleno de XMLs, con estilos, múltiples hojas y metadatos.

	A	B	C	D	E	F	G	
1	España							
2	Actualizado: jueves, 03 enero 2019							
3	Fecha: miércoles, 02 enero 2019							
4								
5	Estación	Provincia	Temperatura máxima (°C)	Temperatura mínima (°C)	Temperatura	Racha (km/h)	Velocidad máxima (km/h)	Pre
6	Estaca de Bares	A Coruña	11.1 (09:30)	9.9 (23:40)	10.5	105 (23:59)	63 (23:00)	
7	As Pontes	A Coruña	11.4 (14:10)	-0.9 (04:30)	5.2			
8	A Coruña	A Coruña	12.3 (13:20)	6.0 (05:40)	9.2	37 (18:20)	24 (17:40)	
9	A Coruña Aeropuerto	A Coruña	12.4 (15:50)	-2.0 (08:50)	5.2	21 (15:30)	13 (07:20)	
10	Carballo, Depuradora	A Coruña	12.5 (15:30)	-2.0 (09:20)	5.3			
11	Cabo Vilán	A Coruña	11.3 (15:10)	7.2 (07:50)	9.3	54 (15:30)	43 (15:30)	
12	Vimianzo	A Coruña	10.8 (15:40)	3.0 (01:00)	6.9			
13	Fisterra	A Coruña	10.9 (15:10)	5.7 (09:20)	8.3	58 (10:10)	44 (05:10)	
14	Mazaricos	A Coruña	10.0 (16:10)	0.8 (04:00)	5.4			
15	Sobrado	A Coruña	10.3 (15:30)	-1.8 (03:40)	4.3			
16	Santiago de Compostela Aeropuerto	A Coruña	11.8 (15:50)	1.3 (04:00)	6.6	18 (13:10)	7 (13:20)	
17	Noia	A Coruña	12.2 (16:10)	4.2 (07:50)	8.2	26 (12:20)	9 (12:40)	
18	Boiro	A Coruña	15.9 (15:40)	0.9 (05:40)	8.4			
19	Padrón	A Coruña	13.8 (15:30)	-0.6 (07:40)	6.6	34 (12:20)	21 (12:20)	
20	Santiago de Compostela	A Coruña	13.1 (16:10)	1.5 (07:30)	7.3	21 (20:10)	9 (01:40)	

Esto hace que la extracción sea **más lenta** y requiera parámetros específicos para navegar dentro del archivo.

pd.read_excel()

Esta es la función que utilizaremos para leer archivos de Excel desde Pandas.

Pero requiere **instalar librerías** para poder hacerlo:

- Para **.xlsx** (Excel actual): openpyxl
- Para **.xls** (Excel 97-2003): xlrd

```
!pip install openpyxl  
!pip install xlrd
```

```
Requirement already satisfied: openpyxl in /opt/conda/lib/python3.11/site-packages (3.1.2)  
Requirement already satisfied: et-xmlfile in /opt/conda/lib/python3.11/site-packages (from  
Requirement already satisfied: xlrd in /opt/conda/lib/python3.11/site-packages (2.0.1)
```

Una característica clave de Excel son las **hojas**, ya que un archivo de Excel es un libro que contiene varias hojas.

El parámetro que podemos utilizar la seleccionar la hoja a leer es **sheet_name**.

Si no indicamos este parámetro por defecto leerá la primera hoja.

```
# Leemos la primera hoja
df = pd.read_excel('./Aemet2019-01-02.xls')

# Leemos la hoja llamada Sevilla
df = pd.read_excel(
    './Aemet2019-01-02.xls',
    sheet_name="Sevilla")

# Leemos la tercera hoja pasando un entero
df = pd.read_excel(
    './Aemet2019-01-02.xls',
    sheet_name=2)
```

Si pasamos el parámetro **sheet_name=None**, leerá todas las hojas, pero devolverá un **diccionario de DataFrames** en lugar de un solo DataFrame.

```
dict_df = pd.read_excel(  
    "Aemet2019-01-02.xls",  
    sheet_name=None  
)  
print(type(dict_df))  
  
<class 'dict'>
```

Normalmente los archivos Excel tienen líneas con logotipos, encabezados,.. que proporcionan información visual pero que no es relevante para extraer datos.

En estos casos podemos usar alguno de los siguientes parámetros:

- **header**: indica la fila en la que se encuentra el encabezado de los datos. Por ejemplo: `header=3`
- **usecols**: indica las columnas que hay que leer. Ejemplo: `usecols='A:F'`
- **skiprows**: salta las primeras filas antes de los datos. Ejemplo: `skiprows=5`
- **skipfooter**: salta las filas del final de la hoja. Ejemplo: `skipfooter=2`

Ejemplo de carga de datos de Excel:

```
df = pd.read_excel(  
    'nomina_rrhh.xlsx',  
    sheet_name='Plantilla_Activa',  
    engine='openpyxl',  
    header=2,  
    usecols='A:E',  
    dtype={'ID_Empleado': str}  
)  
print(df.head())
```

Cargamos solo una hoja

Indicamos explícitamente el motor. Opcional, pero recomendado

Los títulos están en la 3ª fila

Solo cargamos las columnas A hasta la E

En este ejemplo, indico que la columna ID_Empleado se leerá como cadena

2

EXTRACCIÓN: CONECTANDO CON LAS FUENTES

2.3

FORMATOS DE ARCHIVOS COMUNES: JSON

El formato **JSON (JavaScript Object Notation)** es el estándar actual para el intercambio de datos en la web (APIs) y bases de datos NoSQL.

A diferencia de CSV o Excel, JSON es **semi-estructurado y jerárquico** (anidado), lo que supone un reto único: hay que *aplanarlo* para que quepa en una tabla.

```
1 {  
2   "count": 7,  
3   "items": ["socks", "pants", "shirts", "hats"],  
4   "manufacturer": {  
5     "name": "Molly's Seamstress Shop",  
6     "id": 39233,  
7     "location": {  
8       "address": "123 Pickleton Dr.",  
9       "city": "Tucson",  
10      "state": "AZ",  
11      "zip": 85705  
12    }  
13  },  
14  "total_price": "$393.23",  
15  "purchase_date": "2022-05-30",  
16  "country": "USA"  
17 }
```

pd.read_json()

Con esta función podemos leer un archivo JSON

```
[  
  {"id": 1, "producto": "Laptop", "precio": 1200},  
  {"id": 2, "producto": "Mouse", "precio": 25}  
]
```

```
# Pandas infiere que es una lista de registros  
df = pd.read_json('datos.json')  
print(df)
```

El problema es cuando encontramos un **JSON anidado**, es decir, tiene datos dentro de datos

```
[  
  {  
    "id": 1,  
    "nombre": "Ana",  
    "contacto": {  
      "email": "ana@test.com",  
      "telefono": "555-0001"  
    }  
  }  
]
```

Si en este caso usamos `read_json()` obtendremos una columna llamada `contacto` que contiene diccionarios enteros, lo que no nos sirve.

Necesitamos **aplanar el JSON**.

pd.json_normalize()

Esta función *explota* las jerarquías

```
import json

with open('usuarios.json') as f:
    data = json.load(f)

df = pd.json_normalize(data)

# RESULTADO VISUAL DEL DATAFRAME:
# id | nombre | contacto.email | contacto.telefono
# 1 | Ana | ana@test.com | 555-0001
```

Primero lo cargamos con la librería estándar json

Y luego lo aplanamos

Observa como nombra las columnas para reflejar la jerarquía del archivo JSON

Ejemplo:

```
df = pd.read_json('./components.json')
df
```

	id	nombre	especificaciones
0	101	Laptop Gamer	{'ram': '16GB', 'disco': '1TB SSD'}
1	102	Ratón Inalámbrico	{'ram': 'N/A', 'disco': 'N/A'}

```
with open('components.json') as f:
    data = json.load(f)

df = pd.json_normalize(data)
df
```

	id	nombre	especificaciones.ram	especificaciones.disco
0	101	Laptop Gamer	16GB	1TB SSD
1	102	Ratón Inalámbrico	N/A	N/A

```
{
  {
    "id": 101,
    "nombre": "Laptop Gamer",
    "especificaciones": {
      "ram": "16GB",
      "disco": "1TB SSD"
    }
  },
  {
    "id": 102,
    "nombre": "Ratón Inalámbrico",
    "especificaciones": {
      "ram": "N/A",
      "disco": "N/A"
    }
  }
}
```

El problema es que, si el JSON tiene más niveles de anidamiento no funcionará como esperamos:

```
with open('students.json') as f:
    data = json.load(f)

df = pd.json_normalize(data)
df
```

```
[
  {
    "ciclo": "DAM",
    "curso": "Primero",
    "estudiantes": [
      {"nombre": "Juan", "nota": 8},
      {"nombre": "Lucia", "nota": 9}
    ]
  },
  {
    "ciclo": "ASIR",
    "curso": "Segundo",
    "estudiantes": [
      {"nombre": "Pedro", "nota": 5}
    ]
  }
]
```

	ciclo	curso	estudiantes
0	DAM	Primero	[{'nombre': 'Juan', 'nota': 8}, {'nombre': 'Lu...
1	ASIR	Segundo	[{'nombre': 'Pedro', 'nota': 5}]

Para evitar esto, deberemos utilizar los parámetros:

- **record_path**: sirve para indicar a Pandas de dónde queremos sacar cada fila de la tabla final

```
df = pd.json_normalize(  
    data,  
    record_path=['estudiantes']  
)  
df
```

	nombre	nota
0	Juan	8
1	Lucia	9
2	Pedro	5

Observa que en este caso cargamos los datos que hay dentro de estudiantes, pero perdemos los datos de ciclo y curso.

- **meta:** con este parámetro indicamos qué datos del nivel superior queremos pegar al lado de cada estudiante.

```
df = pd.json_normalize(  
    data,  
    record_path=['estudiantes'],  
    meta=['ciclo', 'curso']  
)  
df
```

	nombre	nota	ciclo	curso
0	Juan	8	DAM	Primero
1	Lucia	9	DAM	Primero
2	Pedro	5	ASIR	Segundo

En este caso conservo los
datos del nodo padre

JSON Lines

Algo habitual es encontrar archivos con el formato **JSON Lines**, donde cada línea del archivo es un objeto JSON independiente en lugar de ser todo el archivo un solo objeto gigante.

```
{ "id": "0", "question": "-5 * -6", "answer": "30", "num_terms": 2, "num_digits": 1 }
{ "id": "1", "question": "30 * 44 =", "answer": "1320", "num_terms": 2, "num_digits": 2 }
{ "id": "2", "question": "621242 - 793732 = ?", "answer": "-172490", "num_terms": 2, "num_digits": 6 }
{ "id": "3", "question": "What is 97 + -41 + -19 - -31 + -45 * 84?", "answer": "-3712", "num_terms": 6, "num_digits": 2 }
{ "id": "4", "question": "( -71 + 75 ) - -1 - 36 =", "answer": "-31", "num_terms": 4, "num_digits": 2 }
{ "id": "5", "question": "- ( - ( - ( -5 - 8 ) * 1 ) ) - -6 - -2 = ?", "answer": "21", "num_terms": 5, "num_digits": 1 }
{ "id": "6", "question": "( 69 - 20 * -86 * -84 ) - -72", "answer": "-144339", "num_terms": 5, "num_digits": 2 }
{ "id": "7", "question": "5 - -45", "answer": "50", "num_terms": 2, "num_digits": 2 }
{ "id": "8", "question": "What is -( -( 330191 - -907543 ) + 617155 )?", "answer": "620579", "num_terms": 3, "num_digits": 6 }
{ "id": "9", "question": "-8 - -3 * -6 =", "answer": "-26", "num_terms": 3, "num_digits": 1 }
{ "id": "10", "question": "What is 89 - 6 * -31?", "answer": "275", "num_terms": 3, "num_digits": 2 }
```

Para hacer más eficiente la carga de datos es necesario utilizar el parámetro `lines=True`.

```
df = pd.read_json(  
    'logs_servidor.json',  
    lines=True)
```

2

EXTRACCIÓN: CONECTANDO CON LAS FUENTES

2.4

FORMATOS DE ARCHIVOS COMUNES: XML

Otro tipo de archivo que encontraremos habitualmente es **XML (Extensible Markup Lenguaje)**.

Al igual que pasa con JSON, XML es **jerárquico** mientras que Pandas es tabular, por lo que necesitaremos **aplanar** esa estructura.

```
<?xml version="1.0"?>
- <birds>
  - <owl id="1201">
    <species>Bubo bubo</species>
    <name>Eagle Owl</name>
    <region>Eurasia</region>
  </owl>
  - <owl id="1202">
    <species>Strix occidentalis</species>
    <name>Spotted Owl</name>
    <region>North America</region>
  </owl>
</birds>
```

Antes de leer un fichero XML deberemos tener instalada la librería `lxml`.

```
!pip install lxml
```

pd.read_xml()

Este es el método más sencillo.

Pandas intentará detectar automáticamente la estructura repetitiva (filas) y los datos (columnas)

```
<biblioteca>
  <libro categoria="ficción">
    <titulo>El Quijote</titulo>
    <autor>Cervantes</autor>
    <precio>20.50</precio>
  </libro>
  <libro categoria="educación">
    <titulo>Python Data Science</titulo>
    <autor>VanderPlas</autor>
    <precio>45.00</precio>
  </libro>
</biblioteca>
```

```
df = pd.read_xml('biblioteca.xml')
```

```
df
```

	categoria	titulo	autor	precio
0	ficcion	El Quijote	Cervantes	20.5
1	educacion	Python Data Science	VanderPlas	45.0

El problema es cuando hay niveles de **anidamiento**.

```
<empresa>
  <departamento nombre="IT">
    <empleados>
      <empleado id="1">Juan</empleado>
      <empleado id="2">Ana</empleado>
    </empleados>
  </departamento>
  <departamento nombre="HR"> ... </departamento>
</empresa>
```

```
df = pd.read_xml('biblioteca.xml')
df
```

	nombre	empleados	departamento
0	IT	\n	None
1	HR	None	...

En esos casos debemos utilizar el parámetro **xpath** para indicar dónde debe buscar Pandas los datos.

```
df = pd.read_xml(  
    'data.xml',  
    xpath="//*[empleado]"  
)  
df
```

Le indicamos que busque la etiqueta llamada *empleado*

	id	empleado
0	1	Juan
1	2	Ana

```
<tienda>
  <seccion nombre="electronica">
    <producto id="A001">
      <nombre>Laptop</nombre>
      <precio moneda="USD">800</precio>
    </producto>
  </seccion>
  <seccion nombre="ropa">
    <producto id="B005">
      <nombre>Camiseta</nombre>
      <precio moneda="EUR">20</precio>
    </producto>
  </seccion>
</tienda>
```

Hay 4 formas de indicar la ruta:

- **Rutas absolutas:** comienza por barra inclinada simple (/) e indica la ruta desde la raíz hasta el elemento

```
df = pd.read_xml(  
    'data.xml',  
    xpath="/tienda/seccion/producto"  
)  
df
```

	id	nombre	precio
0	A001	Laptop	800
1	B005	Camiseta	20

- **Rutas relativas:** empieza con doble balla (//) y le dice al sistema: *busca este elemento donde sea que esté, sin importar la profundidad.*

```
df = pd.read_xml(  
    'data.xml',  
    xpath="//producto"  
)  
df
```

	id	nombre	precio
0	A001	Laptop	800
1	B005	Camiseta	20

- **Predicados:** los corchetes funcionan como un WHERE de AQL, permiten seleccionar nodos basándonos en condiciones específicas. Se indican por corchetes. Ejemplo:

```
df = pd.read_xml(  
    'data.xml',  
    xpath='//seccion[@nombre="electronica"]/producto'  
)  
df
```

	id	nombre	precio
0	A001	Laptop	800

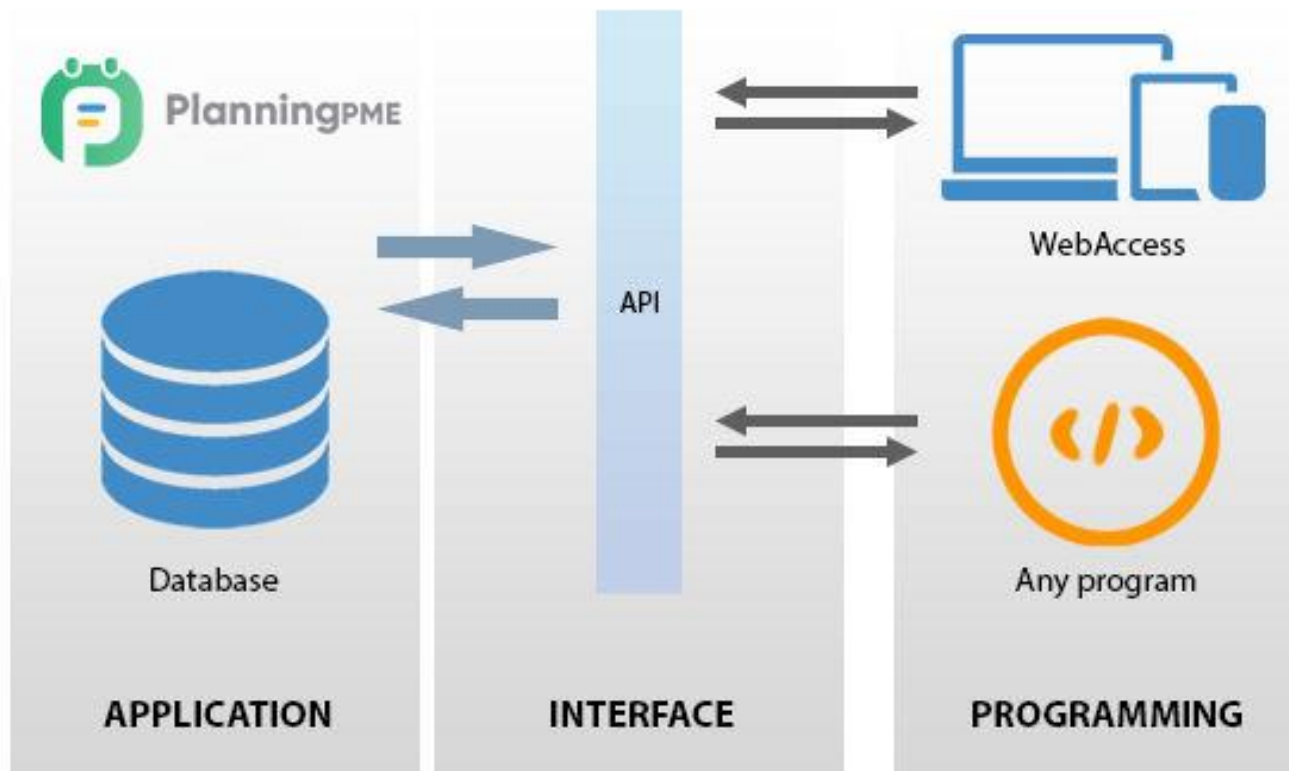
2

EXTRACCIÓN: CONECTANDO CON LAS FUENTES

2.5

LECTURA DE DATOS DE APIs REST

Una **API** (***Application Programming Interface***) es un conjunto de funcionalidades o recursos que expone un sistema para poder interactuar con él desde otro sistema, independientemente del lenguaje de programación o tecnología de cada uno de ellos.



REST (Representational State Transfer) es un estilo de arquitectura del software para comunicaciones cliente servidor apoyado en el **protocolo HTTP**.

REST se basa en tres conceptos clave:

- URLs
- Métodos HTTP
- Estados de respuesta.

URL (Uniform Resource Locator)

Una **URL** es la dirección que se le da a un recurso en la red. REST redefine este concepto utilizándolo para identificar recursos, pero también asignándoles nombres representativos.

Así, las consultas a la API son fácilmente comprensibles.

Por ejemplo:

- <https://swapi.dev/api/people/1/>
- https://api.twitter.com/2/users/:id/timelines/reverse_chronological

Métodos HTTP

Los métodos HTTP se utilizan para indicar qué se quiere hacer con un recurso determinado.

Se utilizan cuatro métodos principalmente, asociados con las operaciones CRUD:

- **GET:** para obtener o leer un recurso.
- **PUT:** actualiza o reemplaza un recurso
- **DELETE:** elimina un recurso del servidor
- **POST:** crea un recurso en el servidor

Estados de respuesta

El resultado de la consulta a la API se indica en el campo de estado de la respuesta HTTP.

Los estados definidos por el estándar HTTP son:

- 1xx Informational
- 2xx Success
- 3xx Redirection
- 4xx Client Error
- 5xx Server Error

Algunos ejemplos:

- **200 (OK)**: la operación solicitada se ha realizado con éxito
- **400 (Bad Request)**: código de error genérico cuando no se adapta ningún otro.
- **401 (Unauthorized)**: el usuario no ha facilitado el método de autenticación requerido por la API y no tiene acceso al recurso.
- **404 (Not Found)**: indica que la API REST no puede mapear la URI con un recurso, pero puede que sí pueda en un futuro, por lo que sí se permitirían futuras solicitudes.
- **500 (Server Error)**: código genérico para indicar algún tipo de error en el servidor.

MORE
INFO



<https://restfulapi.net/http-status-codes/>

No existe función análoga a `pd.read_csv()` que lea datos de una API.

La carga de datos de una API requiere dos pasos:

- **Extracción:** usamos la librería `requests` para obtener los datos (usualmente en formato JSON).
- **Transformación:** usamos Pandas para convertir el JSON en un `DataFrame`

La librería requests

Esta librería es el estándar de Python para interactuar con la web.

Lo primero que hay que hacer es instalarla si no la tenemos en nuestro equipo.

```
!pip install requests
```

requests.get()

Esta función nos permitirá realizar solicitudes GET a un recurso web.

```
import requests
```

```
url = 'https://swapi.info/api/people'
```

```
response = requests.get(url)
```

No obtenemos los datos, sino la respuesta.

```
if response.status_code == 200:  
    print("¡Éxito! Conexión establecida.")  
    datos = response.json()  
    print(datos)
```

La propiedad `status_code` es fundamental para saber qué tipo de respuesta tenemos

```
else:  
    print(f"Error: {response.status_code}")
```

Si el estado es 200 podemos leer los datos con `json()` o `text()`

```
¡Éxito! Conexión establecida.
```

```
[{'name': 'Luke Skywalker', 'height': '172', 'mass': '77', 'hair_cc
```

Y ya podríamos usar las funciones de JSON para importarlo a Pandas

```
pd.json_normalize(datos)
```

	name	height	mass	hair_color	skin_color	eye_color	birth_year	gender	homeworld	films	species	vehicles	starships
0	Luke Skywalker	172	77	blond	fair	blue	19BBY	male	https://swapi.info/api/planets/1	https://swapi.info/api/films/1 , https://swapi...	https://swapi.info/api/species/1	https://swapi.info/api/vehicles/14 , https://s...	https://swapi.info/api/starships/12 , https://...
1	C-3PO	167	75	n/a	gold	yellow	112BBY	n/a	https://swapi.info/api/planets/1	https://swapi.info/api/films/1 , https://swapi...	https://swapi.info/api/species/2	https://swapi...	https://swapi...

Algo importante al consumir datos de una API REST es el **manejo de tiempos de espera**, algo crucial para evitar que el script se quede colgado si el servidor remoto acepta la conexión, pero nunca envía datos.

Esto hará que el script se quede esperando infinitamente.

Por ello, es conveniente usar el parámetro **timeout**, que recibe una tupla de la forma `timeout = (connect, read)`.

- Connect timeout: tiempo máximo para establecer la conexión con el servidor.
- Read timeout: tiempo máximo de espera entre bytes enviados por el servidor.

Tengo que importar las excepciones

Establezco tiempos de espera de 3 seg (*conexión*) y 10 seg (*lectura datos*)

```
import requests
from requests.exceptions import Timeout, RequestException

url = "https://swapi.info/api/people"
try:
    response = requests.get(url, timeout=(3, 10))
    response.raise_for_status()
    print("Datos recibidos correctamente")
except Timeout:
    print("Error: El servidor tardó demasiado en responder")
except RequestException as e:
    # Cualquier otro error
    print(f"Error en la petición: {e}")
```

`raise_for_status()` es un método de response que verifica el código de estado HTTP de la respuesta y, si indica error, lanza una excepción

Esta excepción la lanza el get por el parámetro timeout

Y esta el `raise_for_satus()`

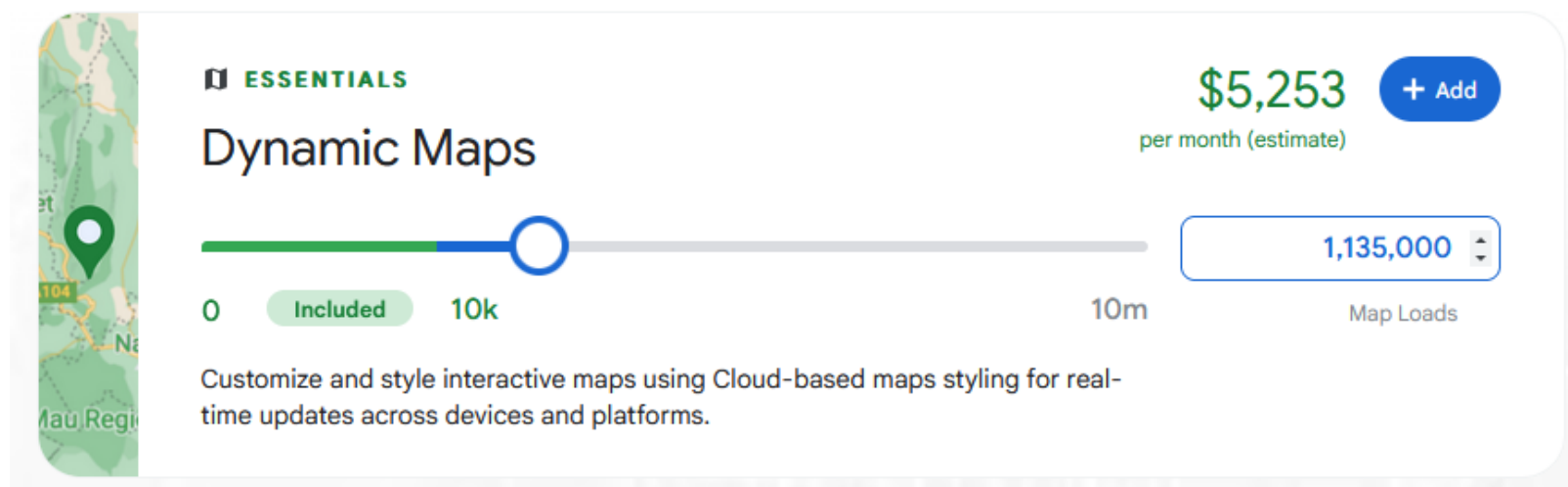
Autenticación

En el caso de una API simple con el uso del método `get()` es suficiente, sin embargo, lo habitual es que las APIs requieran algún tipo de autenticación, que puede ser:

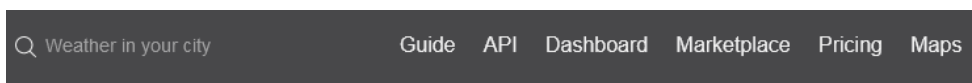
- Una **API Key** que hay que enviar en la *Query String*
- Un **Bearer Token**, que se debe enviar en la cabecera
- **Autenticación Básica**

API Key

Muchos servicios emiten una **API Key**, que básicamente es una secuencia numérica de caracteres, que **vincula nuestras solicitudes a la API con nuestra cuenta personal**, por ejemplo, para contabilizar el número de peticiones y limitarlas o cobrarlas.



Aunque cada API especifica cómo se debe enviar la API Key, es habitual enviarla como **parámetro en la URI** en el denominado *Query String*.



Current & Forecast weather data collection

Current Weather Data

[API doc](#)[Subscribe](#)

- Access current weather data for any location including over 200,000 cities
- We collect and process weather data from different sources such as global and local weather models, satellites, radars and a vast network of weather stations
- JSON, XML and HTML formats

Hourly Forecast 4 days

[API doc](#)[Subscribe](#)

- Hourly forecast is available for 4 days
- Forecast weather data for 96 timestamps
- JSON and XML formats
- Included in the Developer, Professional and Enterprise subscription plans

API call

```
https://api.openweathermap.org/data/2.5/weather?lat={lat}&lon={lon}&appid={API key}
```



Podemos **probar las consultas GET** escribiéndolas directamente en el navegador, aunque lo habitual si queremos hacer muchas pruebas será usar una aplicación específica como **Postman** o extensiones que ha disponibles para todos los navegadores.

The screenshot displays a web browser window with the URL `api.openweathermap.org/data/2.5/weather?lat=42.60&lon=5.57&appid=ae6ba5815beaa93f212d18f38d6f5c69`. The `appid` parameter is highlighted with a red box. Below the URL bar, the JSON response is shown in a structured format:

```
coord:
  lon: 5.57
  lat: 42.6
weather:
  0:
    id: 804
    main: "Clouds"
    description: "overcast clouds"
    icon: "04n"
  base: "stations"
  main:
    temp: 285.87
```

Overlaid on the browser is the Postman interface. The 'Fintech Banking' collection is selected, and the 'GET Overview' endpoint is active. The request is a GET call to `{{base_url}}/accounts/:accountNumber`. The response is a 200 OK status with a body containing JSON data:

```
{
  "balance": {
    "available": 572.18,
    "present": 629.95
  },
  "details": {
    ...
  }
}
```

Para incluir los parámetros de la *Query string* en Pandas debemos hacerlo mediante un diccionario usando el parámetro `params`.

```
params = {  
    "appid": "ae6ba[REDACTED]",  
    "lat": 42.60,  
    "lon": 5.57  
}  
  
url = "https://api.openweathermap.org/data/2.5/weather"  
  
response = requests.get(url, params=params)  
response.json()
```

```
{'coord': {'lon': 5.57, 'lat': 42.6},  
 'weather': [{'id': 804,  
   'main': 'Clouds',  
   'description': 'overcast clouds',  
   'icon': '04n'}],  
 'base': 'stations',  
 'main': {'temp': 285.87,  
   'feels_like': 285.77,  
   'temp_min': 285.77,  
   'temp_max': 285.77,  
   'pressure': 1013.0,  
   'humidity': 77,  
   'wind_speed': 3.6,  
   'wind_deg': 140,  
   'clouds': 100},  
 'visibility': 1000,  
 'sys': {'type': 'weather',  
   'id': 1,  
   'country': 'es',  
   'sunrise': 1181118.5,  
   'sunset': 1181181.5},  
 'timezone': 3600}
```

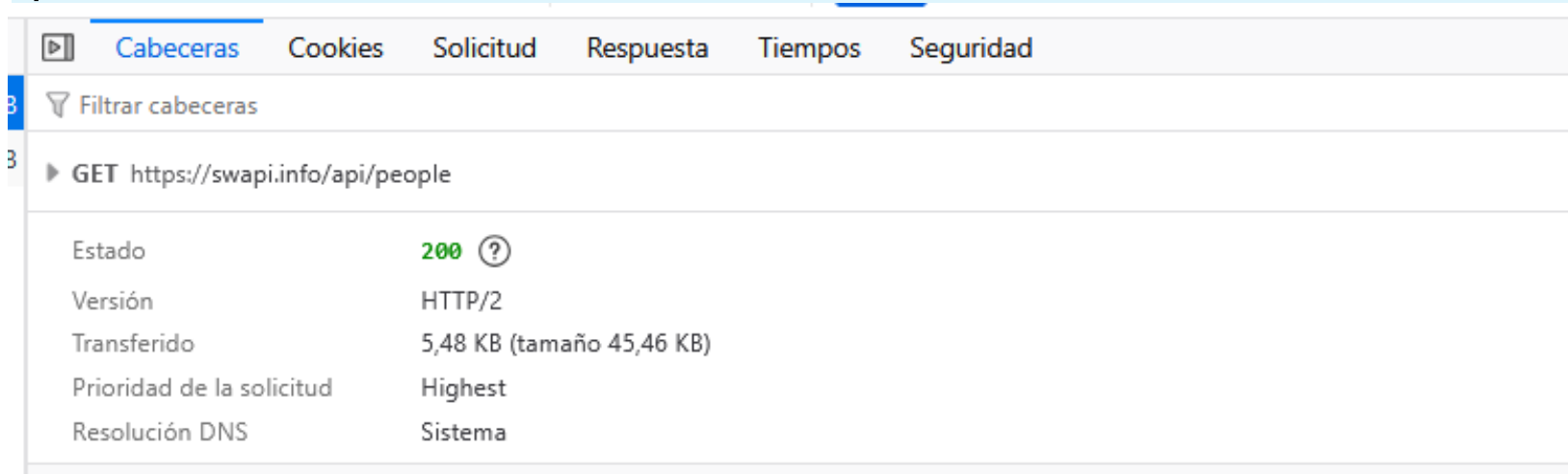
Bearer Token

El **Bearer token** es una alternativa a la API Key que se envía en el **header** del mensaje.

Características:

- Se suele usar en **OAuth**
- Al contrario que las contraseñas, los tokens suelen tener una **fecha de expiración**.
- Es **sin estado**, lo que significa que no necesita recordar que has hecho login. El servidor simplemente lee el token, verifica la firma criptográfica y te deja pasar. Esto hace que las APIs sean rápidas.

El token se envía en el header o las **cabeceras HTTP**, que son **pares nombre:valor** que se intercambian entre el navegador y el servidor en cada petición.

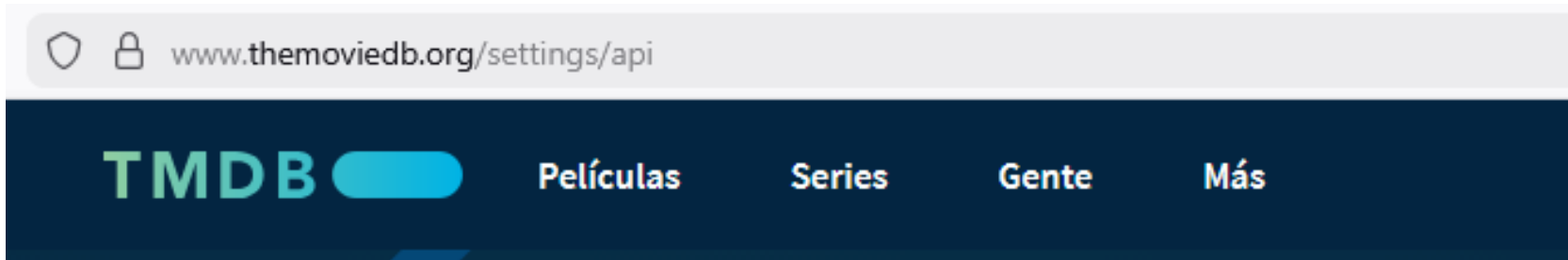


Algunas de las claves que podemos enviar en las cabeceras son:

- **Authorization:** con el JWT, Bearer Tokens o credenciales básicas
- **Content-Type:** define el formato del cuerpo (p.e. application/json)
- **User-Agent:** con el tipo de navegador
- **Accept:** fuerza al servidor a contestar en un formato determinado

Las cabeceras se pueden personalizar en las consultas con requests usando el parámetro **headers**.

Veamos un ejemplo con The Movie Database API, que usa este tipo de autenticación.



Si miramos la ayuda vemos que el token debe enviarse en el header Authorization precedido de la palabra Bearer

cURL Example

```
curl --request GET \  
  --url 'https://api.themoviedb.org/3/movie/11' \  
  --header 'Authorization: Bearer <<access_token>>'
```

La carga de los datos sería así:

```
BEARER_TOKEN='eyJhbGciOiJIUzI1NiJ9.eyJhdWQiOiJkNDI0ZTcwZjgwNDk4YmU0ODU0OGIyNWU2ZmE4NjlmCIIsIm5iZiI
```

```
url="https://api.themoviedb.org/3/movie/popular"
headers={
    "Authorization": f"Bearer {BEARER_TOKEN}",
    "accept": "application/json"
}

response = requests.get(url, headers=headers)

data = response.json()
data
```

```
{'page': 1,
 'results': [{'adult': False,
 'backdrop_path': '/3F2EXWF1thX0BdrVaKvnm6mAhqh.jpg',
 'genre_ids': [28, 53, 80],
 'id': 1306368,
 'original_language': 'en',
 'original_title': 'The Rip',
 'overview': 'Trust frays when a team of Miami cops discovers millions in cash inside a run-down stash house, calling everyone – and everything – into question.',
 ...}]}
```

Basic Auth

Este método es habitual en sistemas antiguos (como JIRA). En estos casos debemos utilizar el parámetro **auth**

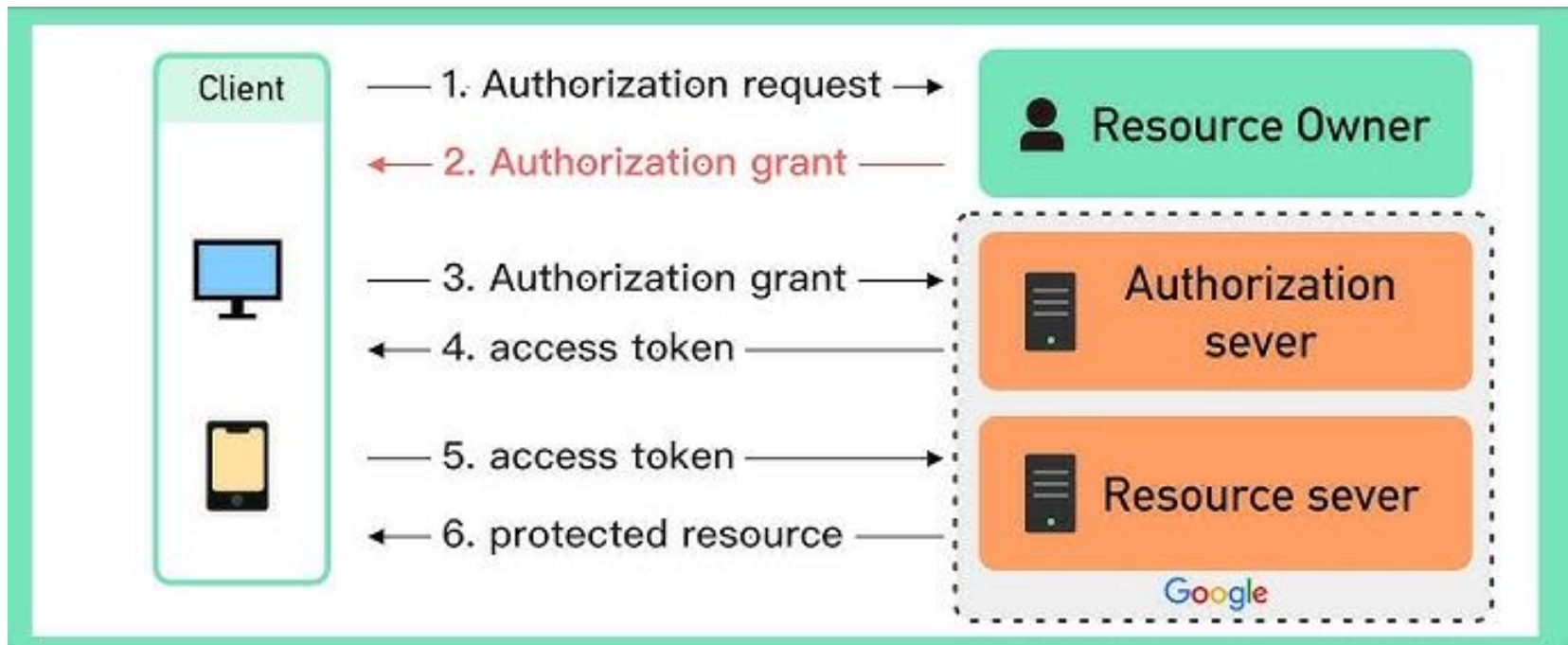
```
from requests.auth import HTTPBasicAuth

user = "mi_usuario"
password = "mi_password"

response = requests.get(url, auth=HTTPBasicAuth(user, password))
```

OAuth

A diferencia de las API Keys simples, que son fijas y eternas, OAuth 2.0 funciona con **tokens de acceso temporales**



El proceso consta de dos pasos claros:

- **Paso 1 (POST):** se envían las credenciales al servidor de autenticación. Este servidor valida las credenciales y devuelve un token de acceso válido por un tiempo determinado.
- **Paso 2 (GET):** se adjunto el token en el encabezado de todas las peticiones siguientes para obtener los datos.

Ejemplo: la API pública de Spotify utiliza este método de autenticación



Request authorization

The first step is to send a POST request to the `/api/token` endpoint of the *Spotify OAuth 2.0 Service* with the following parameters encoded in `application/x-www-form-urlencoded`:

Body Parameters	Relevance	Value
<code>grant_type</code>	<i>Required</i>	Set it to <code>client_credentials</code> .

The headers of the request must contain the following parameters:

Header Parameter	Relevance	Value
Authorization	<i>Required</i>	Base 64 encoded string that contains the client ID and client secret key. The field must have the format: Authorization: Basic <base64 encoded client_id:client_secret>
Content-Type	<i>Required</i>	Set to <code>application/x-www-form-urlencoded</code> .

<https://developer.spotify.com/documentation/web-api/tutorials/client-credentials-flow>

```
auth_str = f"{CLIENT_ID}:{CLIENT_SECRET}"
auth_b64 = base64.b64encode(auth_str.encode()).decode()

headers = {
    'Authorization': f'Basic {auth_b64}',
    'Content-Type': 'application/x-www-form-urlencoded'
}

data = {'grant_type': 'client_credentials'}

response = requests.post(AUTH_URL, headers=headers, data=data)
response.raise_for_status()

# La respuesta incluye 'access_token', 'token_type', 'expires_in' (segundos)
token_info = response.json()
return token_info['access_token']
```


1. `encode()` convierte la cadena en una secuencia de bytes

```
auth_str = f"{CLIENT_ID}:{CLIENT_SECRET}"  
auth_b64 = base64.b64encode(auth_str.encode()).decode()
```

```
headers = {  
    'Authorization': f'Basic {auth_b64}',  
    'Content-Type': 'application/x-www-form-urlencoded'  
}
```

```
data =
```

```
response = requests.post(AUTH_URL, headers=headers, data=data)  
response.raise_for_status()
```

```
# La respuesta incluye 'access token', 'token type', 'expires in' (segundos)
```

```
token_info = response.json()  
return token_info
```

3. `decode()` vuelve a convertir la secuencia de bytes en una cadena para enviarla

2. `base64.encode()` codifica la secuencia de bytes en base64

Diferencia entre cadenas y cadenas de bytes:

- **(str)**: texto abstracto. Usa estándar Unicode, pero a Python no le interesa cómo se guarda en memoria.
- **(bytes)**: datos crudos. Secuencia inmutable de números enteros entre 0 y 255. Para la máquina, no son letras, son valores binarios

```
auth_str = f"{CLIENT_ID}:{CLIENT_SECRET}"  
auth_b64 = base64.b64encode(auth_str.encode('utf-8'))
```

```
headers = {  
    'Authorization': f'Basic {auth_b64}',  
    'Content-Type': 'application/x-www-form-urlencoded'  
}
```

```
data = {'grant_type': 'client_credentials'}
```

```
response = requests.post(AUTH_URL, headers=headers, data=data)  
response.raise_for_status()
```

```
# La respuesta incluye 'access_token'  
token_info = response.json()  
return token_info['access_token']
```

Parameter	Relevance	Value
Authorization	Required	Base 64 encoded string that contains the client ID and client secret. Authorization: Basic <base64 encoded client_id:client_secret>
Content-Type	Required	Set to application/x-www-form-urlencoded.

Body Parameters	Relevance	Value
grant_type	Required	Set it to client_credentials.

```
auth_str = f"{CLIENT_ID}:{CLIENT_SECRET}"  
auth_b64 = base64.b64encode(auth_str.encode()).decode()
```

```
headers = {  
    'Authorization': f'Basic {auth_b64}',  
    'Content-Type': 'application/x-www-form-urlencoded'  
}
```

Realizo una petición POST para solicitar el token

```
data = {'grant_type': 'password', 'username': '...', 'password': '...'}
```

```
response = requests.post(AUTH_URL, headers=headers, data=data)  
response.raise_for_status()
```

```
# La respuesta incluye 'access_token', 'token_type', 'expires_in' (segundos)  
token_info = response.json()  
return token_info['access_token']
```

Envío las cabeceras y los datos

Y en el mensaje devuelto recojo el token