

# doc-pract-1

April 7, 2022

## 1 Tipología y ciclo de vida de los datos: PRÁCTICA I

Componentes de la práctica:

**Jonás Medina Brito** (jmedinabrit@uoc.edu)

**Miguel Rafael Esteban Martín** (mestebanmart@uoc.edu)

### 1.1 Índice de contenidos

- Contexto
- Título
- Descripción del dataset
- Representación gráfica
- Contenido
- Agradecimientos
- Inspiración
- Licencia
- Código
- Dataset

---

### 1.2 Contexto

---

### 1.3 Título

---

### 1.4 Descripción del dataset

---

### 1.5 Representación gráfica

---

## 1.6 Contenido

### 1.6.1 Recogida de datos

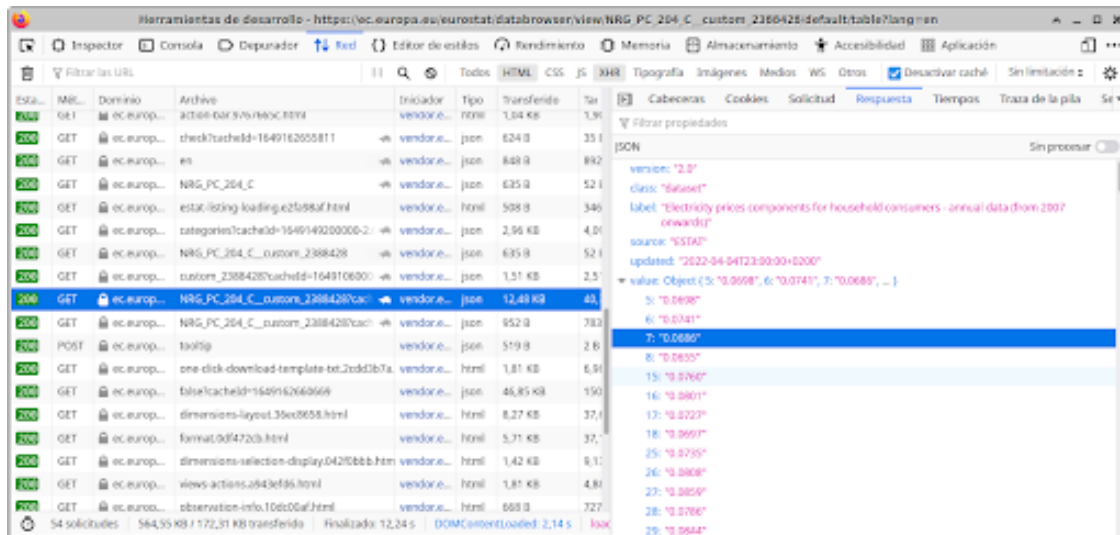
Aunque Eurostat tiene una API pública (1) que permite obtener los diferentes dataset, en esta práctica hemos decido hacer web scraping de los datos que se ofrecen en su “*Data browser*” accesible desde un navegador.

- (1) Eurostat web services: <https://ec.europa.eu/eurostat/en/web/main/data/web-services> (Última vista abril del 2020)

Para explicar el estudio que se ha hecho para la obtención de datos se tomará como ejemplo el dataset de “*Componentes del precio de la electricidad para consumidores domésticos*” ([Electricity prices components for household consumers - annual data \(from 2007 onwards\)](#))



cual es la llamada que devuelve los datos que se representan en la tabla dinámica. Un ejemplo es la petición que es utilizada en la construcción de la tabla con los datos que se muestra en la página web.



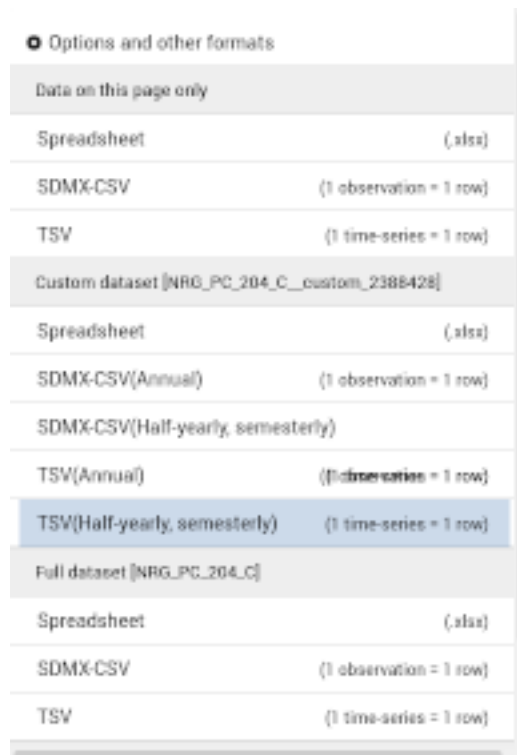
### *Captura de pantalla de la consola del navegador con la petición de datos de dataset y respuesta JSON*

Si se exporta la petición esta es la llamada que se realiza (comando CURL):

```
$ curl 'https://ec.europa.eu/eurostat/databrowser-backend/api/extraction/1.0/LIVE/false/json/en/NRG_PC_204_C_custom_2388428?cacheId=1649106000000-2.6.3%2520-%25202022-03-17%252005%253A43' ...
```

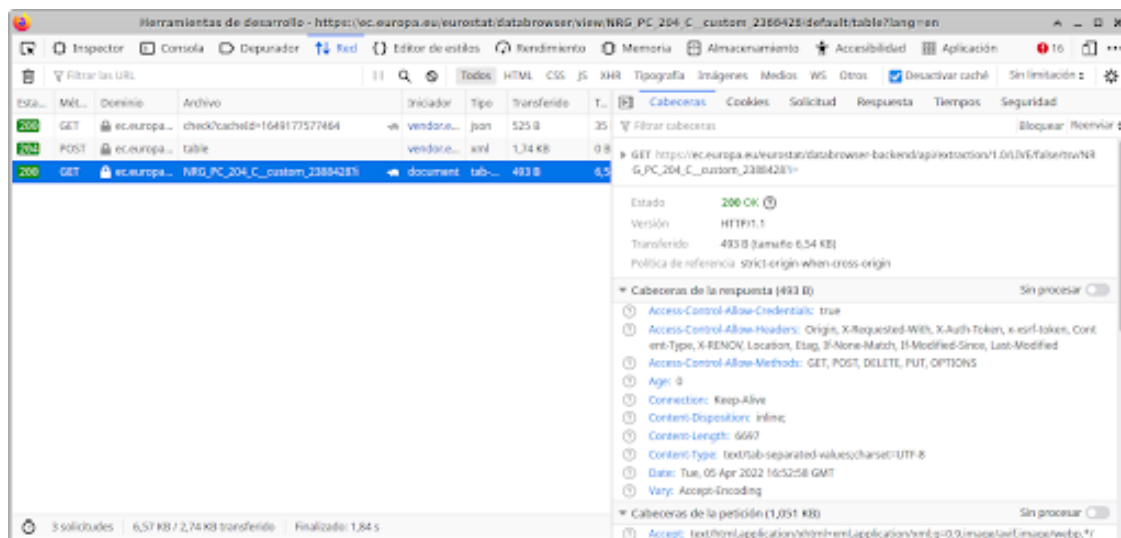
La respuesta de esta llamada es en formato JSON con una estructura específica, difícil de procesar y aprovechar para nuestro propósito de acceder a los datos.

Sin embargo la página tiene un formulario de descarga que permite descargar los dataset, en formato TSV mediante el navegador.



### *Captura de pantalla del enlace de descarga del dataset*

Si se hace click en ese enlace la petición que hace el navegador es la siguiente:



### *Captura de pantalla de la interfaz “Data Browser” de la descarga del informe de los componentes del precio de la electricidad para consumidores domésticos en formato TSV*

Qué si se exporta la petición en formato CURL es:

```
$ curl 'https://ec.europa.eu/eurostat/databrowser-backend/api/extraction/1.0/LIVE/
```

false/tsv/NRG\_PC\_204\_C\_\_custom\_2388428?i' ...

Y que si se ejecuta desde una consola se obtiene los datos en el formato TSV

freq,nrg_cons,nrg_prc,currency,geo\TIME_PERIOD	2012-S2	2013-S2	2014-S2	2015-S2	
A,KWH2500-4999,NETC,EUR,AL	: : : : :	0.0000	0.0000	0.0000	: 0.0000
A,KWH2500-4999,NETC,EUR,AT	: : : : :	0.0606	0.0626	0.0645	0.0639 0.0676
A,KWH2500-4999,NETC,EUR,BA	: : : : :	0.0381	0.0388	0.0367	0.0370 :
A,KWH2500-4999,NETC,EUR,BE	: : : : :	0.1055	0.1116	0.1092	0.1049 0.1046
A,KWH2500-4999,NETC,EUR,BG	: : : : :	0.0232	0.0242	0.0256	0.0265 0.0273
A,KWH2500-4999,NETC,EUR,CY	: : : : :	0.0313	d 0.0319	0.0320	0.0296 0.0272
A,KWH2500-4999,NETC,EUR,CZ	: : : : :	0.0483	0.0520	0.0557	0.0534 0.0425
..					

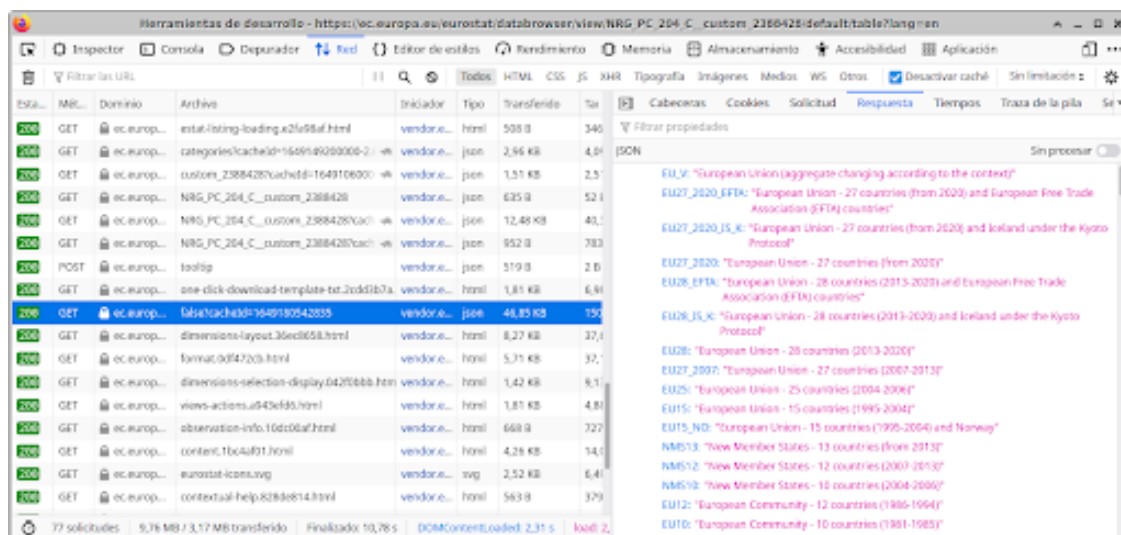
**Llamada HTTP para la obtención del los datos de los dataset** Analizando la petición anterior, la llamada que se hará en el caso práctico para la obtención de los diferentes dataset que componen el dataset principal, será una petición GET HTTP a la URL

[https://ec.europa.eu/eurostat/databrowser-backend/api/extraction/1.0/LIVE/false/tsv/<key\\_dataset>?i](https://ec.europa.eu/eurostat/databrowser-backend/api/extraction/1.0/LIVE/false/tsv/<key_dataset>?i)

Donde:

- **<key\_dataset>** indica el identificador del dataset del sistema de eurostat, identificado como “*online data code*” en cada uno de los conjuntos de datos que son accesibles desde el “*Data Browser*” de Eurostat. Esa llamada es fácilmente procesable desde las utilidades que ofrece la librería PANDAs para la carga de dataset.

**Llamada HTTP para la obtención de los maestros de los países** En los conjuntos de datos descargables, los países están establecidos por un identificador. Para obtener el nombre de los países a partir de su identificador, se utiliza la siguiente petición, también obtenida por ingeniería inversa gracias a las herramientas de desarrollo del navegador.



*Captura de pantalla de la consola del navegador con la petición del maestro*

### *de datos de los países con respuesta en formato JSON*

La llamada HTTP en formato CURL es:

```
$ curl 'https://ec.europa.eu/eurostat/databrowser-backend/api/codelist/LIVE/GEO/getCodeListJson/9.0/ESTAT/en/false?cacheId=1649180542835'...
```

Como respuesta se obtiene un JSON en cuya estructura en la clave `category/label` donde se encuentra el mapeo entre el identificador y el nombre del país.

```
{
  "version": "2.0",
  "class": "dimension",
  "source": "ESTAT",
  "category": {
    "label": {
      "EUR": "Europe",
      "EU": "European Union (EU6-1958, EU9-1973, EU10-1981, EU12-1986, EU15-1995, EU25-2004)",
      "EU_V": "European Union (aggregate changing according to the context)",
      "EU27_2020_EFTA": "European Union - 27 countries (from 2020) and European Free Trade Association",
      "EU27_2020_IS_K": "European Union - 27 countries (from 2020) and Iceland under the EFTA",
      "EU27_2020": "European Union - 27 countries (from 2020)",
      "EU28_EFTA": "European Union - 28 countries (2013-2020) and European Free Trade Association",
      "EU28_IS_K": "European Union - 28 countries (2013-2020) and Iceland under the Kyoto Protocol",
      "EU28": "European Union - 28 countries (2013-2020)",
      "EU27_2007": "European Union - 27 countries (2007-2013)",
      "EU25": "European Union - 25 countries (2004-2006)",
      "EU15": "European Union - 15 countries (1995-2004)",
      ...
    }
  }
}
```

Esta información se obtiene desde el notebook de jupyter gracias a las librerías `requesty json`

#### **1.6.2 Dataset auxiliares obtenidos para la generación del dataset principal**

Los dataset relacionados con los precios del gas y la electricidad en los países europeos que se han utilizado para crear el dataset principal son:

- **Precios del gas (Euro/kWh) para consumidores domésticos.**([Gas prices components for household consumers - annual data](#)). Código del dataset (*online data code*): **NRG\_PC\_202\_C**
- **Precios del gas (Euro/kWh) para empresas.**([Gas prices components for non-household consumers - annual data](#)). Código del dataset (*online data code*): **NRG\_PC\_203\_C**
- **Precio de la electricidad (Euro/kWh) para consumidores domésticos para la banda de consumo entre 2.500 a 4.999 kWh** ([Electricity prices components for household consumers - annual data \(from 2007 onwards\)](#)). Código del dataset (*online data code*): **NRG\_PC\_204\_C**
- **Precio de la electricidad (Euro/kWh) para empresas para la banda de consumo de menos de 20 MWh.**([Electricity prices components for non-household consumers - annual data \(from 2007 onwards\)](#)). Código del dataset (*online data code*): **NRG\_PC\_205\_C**

### 1.6.3 Procesamiento de los dataset auxiliares

Es necesario realizar un proceso de procedimiento de procesamiento y tratamiento de la información en bruto de los dataset. Todos cumplen la misma estructura por lo que el flujo de tratamiento es el mismo, con los siguientes pasos:

- 1) Filtrar los datos. El dataset contiene información adicional que no es necesaria para el alcance de este proyecto. Los filtros comunes a todos los dataset son:
  - Datos anuales
  - Componentes del precio de la energía: *“Energía y suministro”*
  - Moneda: Euro (€)
  - Unidad de medida: kWh
- 2) Obtener el identificador del país. Este dato viene incluido con más información y es necesario extraerla
- 3) Eliminar los espacios de los nombres de las columnas de los dataframe.
- 4) Procesar todas aquellas columnas relativas a los años para que contengan datos numéricos. Algunos datos vienen marcados con flags, como dato confidencial y estimado. Se ha tomado la decisión de que los datos confidenciales se tratan como vacíos y los estimados como valores reales.
- 5) Se añade la columna con la descripción del país.

**Filtrado para los precios del gas (Euro/kWh)** Para ambos dataset relacionados con el gas, se filtrarán los datos que cumplan los filtros comunes que se han indicado más arriba y además:

- Consumo de la energía: En Giga Julios en todas las bandas

**Filtrado para los precios de la electricidad (Euro/kWh) para consumidores domésticos** Se utilizará los filtros comunes y se tendrá en cuenta sólo:

- Consumo de la energía: Consumo entre 2500 kWh y 4999 kWh

**Filtrado para los precios de la electricidad (Euro/kWh) para empresas** Se utilizará los filtros comunes y se tendrá en cuenta sólo:

- Consumo de la energía: Consumo menos de 20 MWh

---

## 1.7 Agradecimientos

---

## 1.8 Inspiración

---

## 1.9 Licencia

---



## 1.10 Código

---

## 1.11 Dataset