

# PRA 1 - Tipología y ciclo de vida de los datos

Maria Dolores Moyano Guerrero y Víctor Cáncer Castillo

## Índice

Contexto . . . . .	2
Título . . . . .	3
Descripción del dataset . . . . .	3
Representación gráfica . . . . .	5
Contenido . . . . .	6
Agradecimientos . . . . .	7
Inspiración . . . . .	7
Licencia . . . . .	8
Código . . . . .	9
Dataset . . . . .	9
Vídeo . . . . .	9
Índice de figuras . . . . .	10

## Contexto

**Explicar en qué contexto se ha recolectado la información. Explicar por qué el sitio web elegido proporciona dicha información.**

Estos datos se han recogido para practicar el *web scraping* en la asignatura de *Tipología y ciclo de vida de los datos* del Máster de ciencia de datos de la UOC.

Cómo (futuros) científicos de datos hemos tenido la curiosidad de estudiar cómo está el mercado laboral actualmente en varias ciudades europeas y americanas. Además hemos querido averiguar en qué lugares el trabajo de científico de datos está más reconocido por las empresas y por lo tanto mejor remunerados.

Para ello hemos obtenido los sueldos que se ofrecen por diferentes empresas utilizando la web <https://glassdoor.es/>, donde los trabajadores pueden informar de su sueldo de manera anónima. Por otro lado hemos extraído datos de la web <https://datosmacro.expansion.com/> donde hay múltiples datos económicos, entre ellos el salario medio, lo cual nos puede mostrar si el trabajo del científico de datos está mejor/peor remunerado que el resto de trabajos en ese país o ciudad.



**Figura 1:** Captura de las páginas webs utilizadas para el estudio

## Título

**Definir un título que sea descriptivo para el dataset.**

El título que describe el dataset final es: comparación de los sueldos para un científico de datos en determinadas ciudades con el sueldo medio de los países donde radican dichas ciudades.

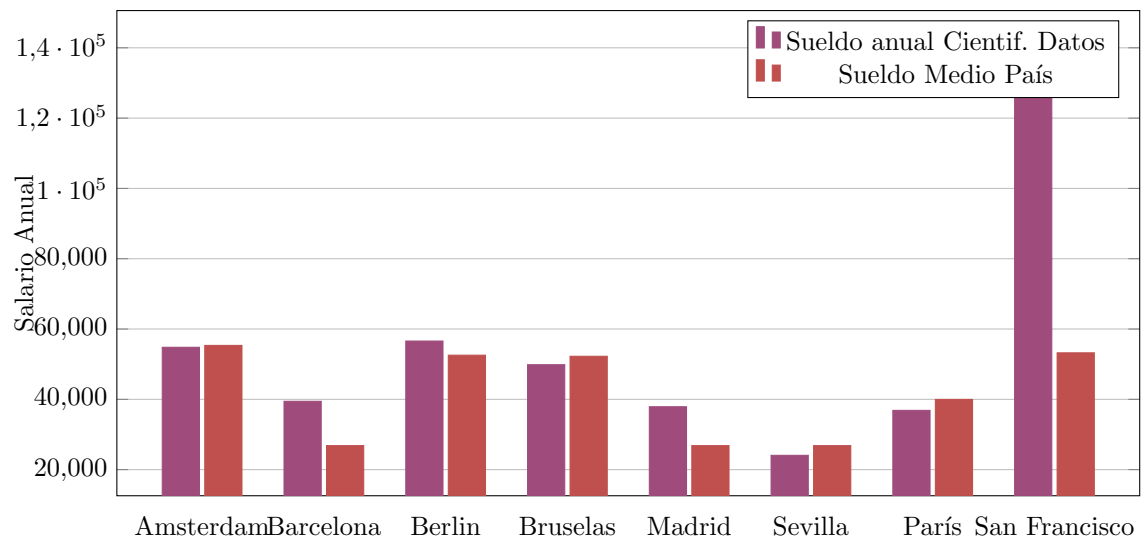
## Descripción del dataset

**Desarrollar una descripción breve del conjunto de datos que se ha extraído. Es necesario que esta descripción tenga sentido con el título elegido.**

Para las ciudades descritas se recogen los datos:

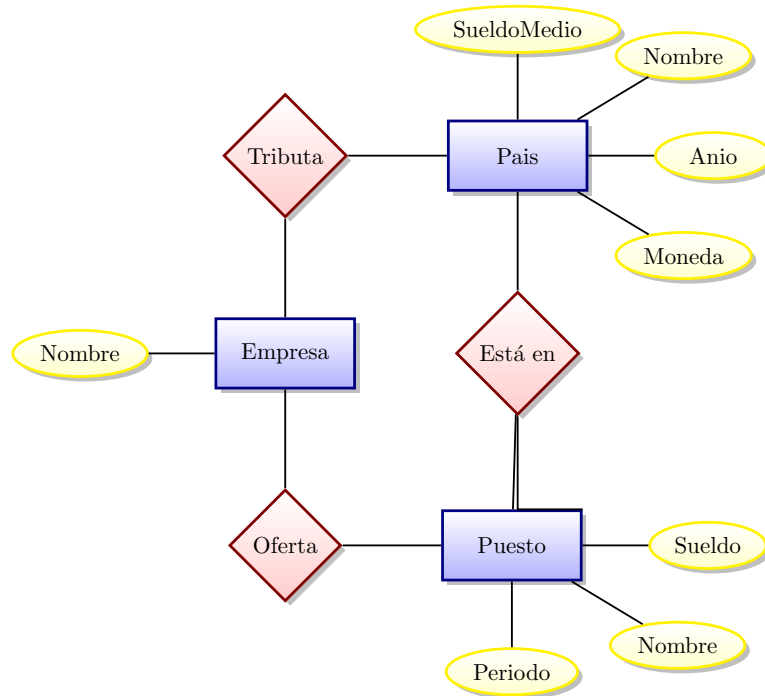
- *Sueldos\_DataScientist\_vs\_SMI*: Comparación de los sueldos de las distintas ciudades con el Sueldo Medio.
- *SMI.csv*: Contiene los sueldos medios de un listado de países.
- *dataset\_sueldos.csv*: Sueldos extraídos para científicos de datos en distintos periodos.
- *dataset\_sueldos\_clean.csv*: Sueldos de un científico de datos para un año.

Los datos obtenidos recogen los datos siguientes:



**Figura 2:** Comparación de los sueldos de científico de datos en las ciudades estudiadas

## Representación gráfica



**Figura 3:** *Diagrama de las entidades principales del proyecto*

Dibujar un esquema o diagrama que identifique el dataset visualmente y el proyecto elegido.

## Contenido

**Explicar los campos que incluye el dataset, el periodo de tiempo de los datos y cómo se han recogido.**

Dataset\_sueldos.csv: se ha obtenido de la página <https://glassdoor.es/> el 6/4/2022. No ha sido necesario recoger los datos en varios días debido a la naturaleza del proyecto.

- **Ciudad:** Ciudad donde está radicado el empleo.
- **Empresa:** Empresa que proporciona el empleo.
- **Sueldo:** Sueldo bruto correspondiente al empleo, ciudad y empresa.
- **Periodo:** Empleo en el periodo indicado, en su caso.

Dataset\_sueldos\_clean.csv: se ha obtenido de la página <https://glassdoor.es/> el 6/4/2022. No ha sido necesario recoger los datos en varios días debido a la naturaleza del proyecto.

- **Ciudad:** Ciudad donde está el empleo.
- **Empresa:** Empresa que proporciona el empleo.
- **Sueldo:** Sueldo bruto correspondiente al empleo, ciudad y empresa.
- **Sueldo Anual:** Sueldo ajustado a un año (algunas empresas, no lo proporcionan anualizado).

SMI.csv: se ha obtenido de la página <https://datosmacro.expansion.com/> el 6/4/2022. No ha sido necesario recoger los datos en varios días debido a la naturaleza del proyecto.

- **País:** País que se corresponde al sueldo medio.
- **Año:** Año de cálculo del sueldo medio.
- **SalMed Local:** Salario medio en moneda local.
- **Moneda:** Moneda del salario local.
- **Salmed \$:** Salario medio en dólares.
- **Salmed €:** Salario medio en euros.
- **Var:** Variación del salario medio.

Sueldos\_DataScientist\_vs\_SMI: Se han calculado a partir de los datos extraídos en las páginas descritas.

- **Ciudad:** Ciudad donde está el empleo.
- **Sueldo anual:** Sueldo anual de los empleos descargados.
- **Suelvo vs Sueldo Medio:** Diferencia entre el sueldo del puesto de Científico de Datos y el sueldo medio del país.

## Agradecimientos

Agradecemos a los propietarios de las páginas <https://glassdoor.es/> y <https://datosmacro.expansion.com/>, por los datos de los que nos hemos podido alimentar en este proyecto de Web Scraping.

## Inspiración

En el actual contexto, donde las comunicaciones permiten buscar trabajo en distintos destinos sin tener que desplazarse, y además es posible trabajar en remoto, los sitios webs que permiten buscar trabajo y comparar sueldos son una herramienta muy útil a la hora de escoger el puesto más conveniente. Los datos recogidos son interesantes porque permiten comparar los sueldos de los científicos de datos en caliente, utilizando el factor de corrección del sueldo medio del país, para poder analizar si efectivamente esos puestos de trabajo son valorados en la ciudad donde se demandan y así tener un mapa claro de la situación, para prestar ayuda a los demandantes de empleo de esta disciplina.

Este proyecto podría ampliarse a la búsqueda de trabajo para otra tipología de puestos, simplemente, modificando la búsqueda.

## Licencia

Se ha escogido la licencia CC BY-SA 4.0 ya que es la más correcta en cuanto a las libertades que se han de cumplir:



### Reconocimiento-Compartirigual 4.0 Internacional (CC BY-SA 4.0)

This is a human-readable summary of (and not a substitute for) the [license](#). [Advertencia](#).

#### Usted es libre de:

**Compartir** — copiar y redistribuir el material en cualquier medio o formato

**Adaptar** — remezclar, transformar y crear a partir del material para cualquier finalidad, incluso comercial.

El licenciador no puede revocar estas libertades mientras cumpla con los términos de la licencia.



#### Bajo las condiciones siguientes:

 **Reconocimiento** — Debe [reconocer adecuadamente](#) la autoría, proporcionar un enlace a la licencia e [indicar si se han realizado cambios](#). Puede hacerlo de cualquier manera razonable, pero no de una manera que sugiera que tiene el apoyo del licenciador o lo recibe por el uso que hace.

 **Compartirigual** — Si remezcla, transforma o crea a partir del material, deberá difundir sus contribuciones bajo la [misma licencia que el original](#).

**No hay restricciones adicionales** — [No puede aplicar](#) términos legales o medidas tecnológicas que legalmente restrinjan realizar aquello que la licencia permite.

#### Avisos:

No tiene que cumplir con la licencia para aquellos elementos del material en el dominio público o cuando su utilización esté permitida por la aplicación de [una excepción o un límite](#).

No se dan garantías. La licencia puede no ofrecer todos los permisos necesarios para la utilización prevista. Por ejemplo, otros derechos como los de [publicidad](#), [privacidad](#), o los [derechos morales](#) pueden limitar el uso del material.



## Código

Para hacer tanto scraping como pre-procesado de datos hemos utilizado Python. El código está disponible en GitHub: [https://github.com/mdmoyano/web\\_scraping/src](https://github.com/mdmoyano/web_scraping/src).

## Dataset

El Dataset está publicado en Zenodo, disponible en [10.5281/zenodo.6408002](https://zenodo.org/record/105281/files/6408002)

## Vídeo

– Link al video de cada uno –

## Índice de figuras

1	Captura de las páginas webs utilizadas para el estudio . . . . .	2
2	Comparación de los sueldos de científico de datos en las ciudades estudiadas . . . . .	4
3	Diagrama de las entidades principales del proyecto . . . . .	5