



JÚLIA  
RODRÍGUEZ

# Proyecto Web Scraping

10 de abril, 2021

# Índice

- 1. Objetivos de negocio**
- 2. Datos que necesitamos**
- 3. Métodos y resultados**
- 4. Obstáculos**
- 5. Aprendizajes**
- 6. Siguietes pasos**

01

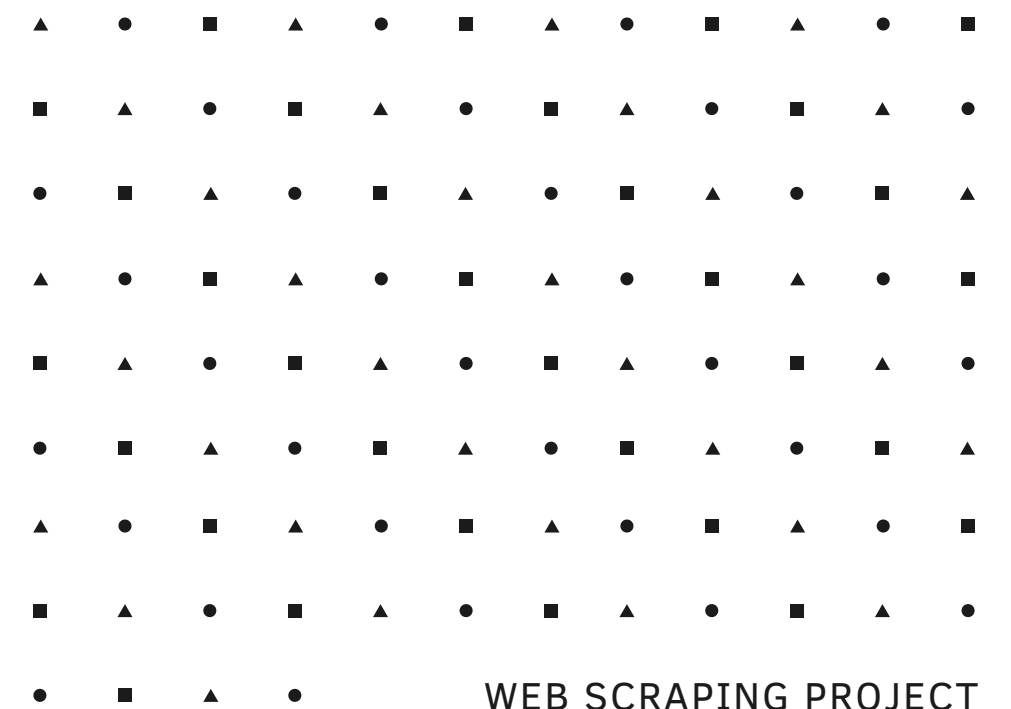
# Objetivos de negocio

## Obtener datos para analizar el posicionamiento orgánico de los productos de Revlon Professional en Amazon España

Objetivo 1: analizar el posicionamiento de los productos cuando alguien busca el nombre de la marca.

Objetivo 2: analizar el posicionamiento de los productos de Revlon Professional y de los principales competidores cuando alguien busca productos de peluquería profesional sin especificar marca.

REVLON  
PROFESSIONAL™



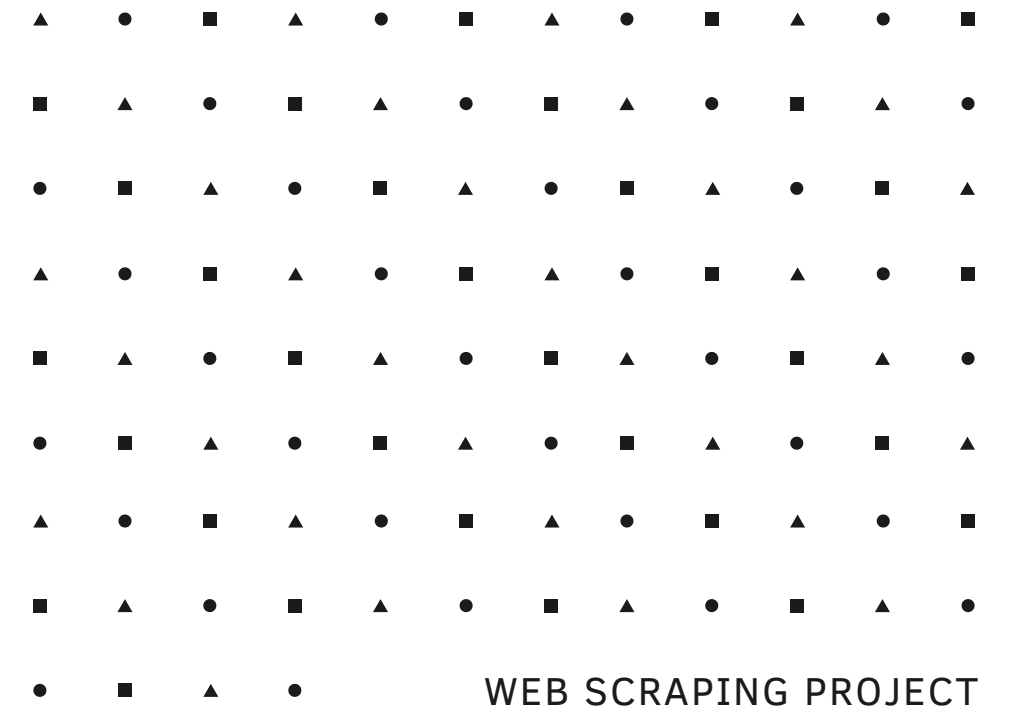
WEB SCRAPING PROJECT

02

# Datos que necesitamos

- Nombre del producto
- Marca
- Precio final
- Precio a granel
- Precio sin descuento
- Rating
- Número de reviews
- Fecha de entrega

REVLON  
PROFESSIONAL™



1-48 de 782 resultados para "revlon professional"

Ordenar por: Destacados ▾

Con derecho a envío gratis

☐ Envío gratis

Envío GRATIS en pedidos superiores a 19€ en Libros o a 29€ en las demás categorías de productos

Día de entrega

☐ Recíbelo mañana

Productos frescos

☐ Amazon Fresh

☐ Supermercado DIA

Departamento

◀ Cualquier departamento

◀ Belleza

◀ Cuidado del cabello

Productos para el cuidado del cabello

Mascarillas de pelo

Champús

Acondicionadores de pelo

Máscaras de tinte de pelo

Valoración media de los clientes

★★★★★ o más

★★★★☆ o más

★★★☆☆ o más

★★☆☆☆ o más

Marca

☐ REVLON PROFESSIONAL

☐ REVLON

☐ EQUAVE

☐ L'Oréal Professionnel

☐ TIGI

☐ TIGI Bed Head

☐ Ckeyin

☐ BED HEAD by TIGI

REVLON

Descubre los productos de Revlon

Descúbrelo REVLO >



Revlon Wonder Woman Super Lustrous Pintalabios (Amazon)

★★★★★ 308



Revlon Wonder Woman Gift Pack

★★★★★ 8



Revlon Wonder Woman Ultra HD Vinyl Mouse Labial (So Shady)

★★★★★ 72

Patrocinado ⓘ



Patrocinado ⓘ

UniqOne Revlon Professional Classico Tratamiento en Spray para Cabello 150 ml y Super10R Mascarilla 300 ml - Pack 300 ml

★★★★★ 71



Revlon Professional UniqOne Champú y Acondicionador 1000 ml

★★★★★ 926

16,99€ (1,70 €/100 ml)

Envío GRATIS



UniqOne Revlon Professional - Tratamiento para el cabello, Coco, 150 ml

★★★★★ 8.810

6,95€ (4,63 €/100 ml) 11,66€

6,60 € con el descuento de Compra



Revlon Professional ProYou Textura de Peinado para Cabello Fino 350 ml

★★★★★ 71

12,99€ (12,99 €/100 g)

12,34 € con el descuento de Compra recurrente

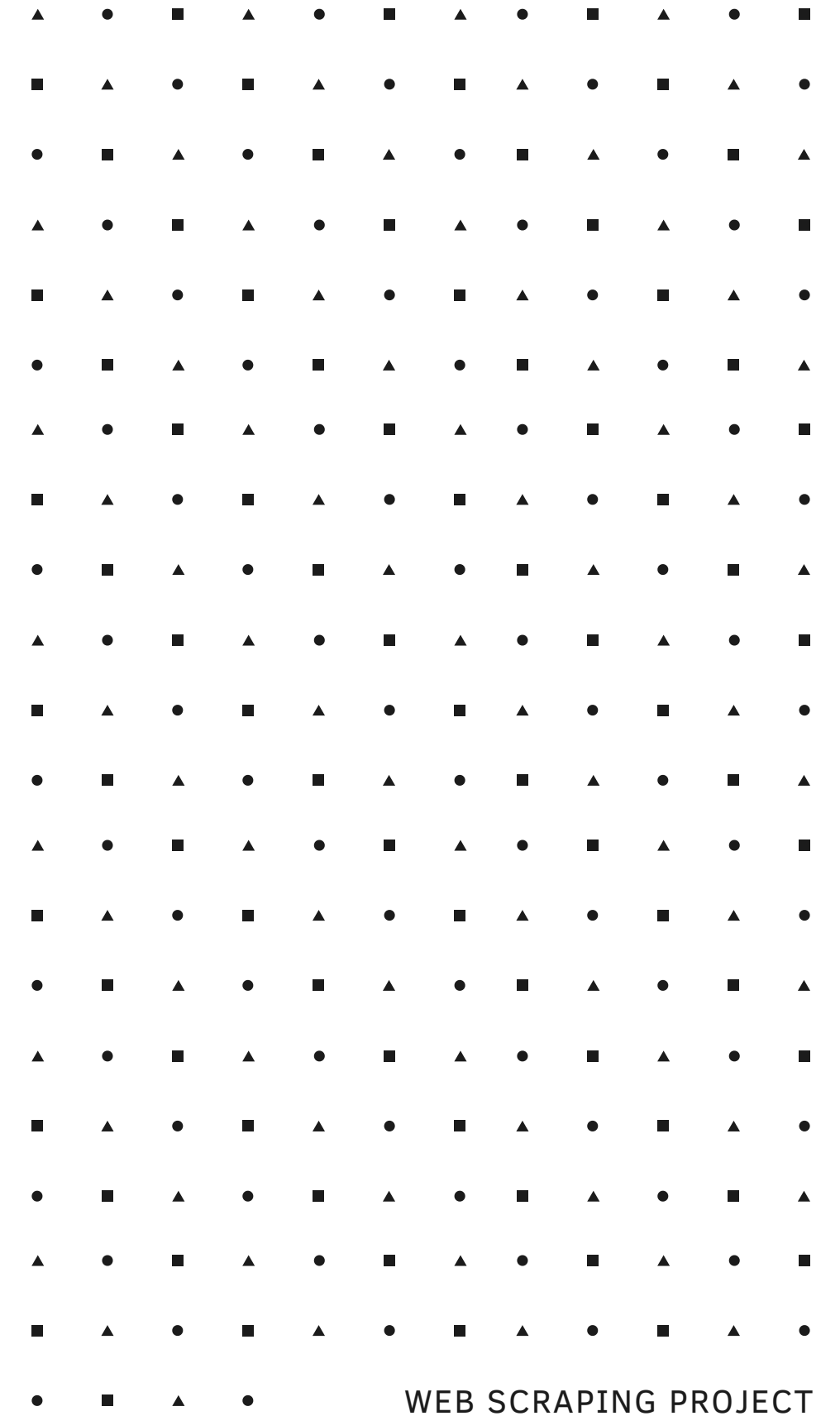




03

# Métodos y resultados

1. BeautifulSoup
2. Scrapy
3. Selenium (descartado)



Para el objetivo 1, necesitamos obtener 796 artículos distribuidos en 17 páginas

1. Obtenemos los XPaths de los productos y de todos los campos que nos interesan.
2. Generamos las URLs de todas las páginas con una f string y un for loop.
3. Parseamos todas las páginas con un for loop convirtiendo cada producto en un diccionario.
4. Lo convertimos todo a DataFrame y lo exportamos como csv.

Resultado:

```
# Checking there are no errors  
df2.head()
```

	prod_name	price	old_price	price_ml	stars	reviews	delivery_date
0	Revlon Professional UniqOne Champú y Acondicio...	NaN	NaN	NaN	4.5	3003.0	None
1	UniqOne Revlon Professional - Tratamiento para...	6.95	11.66	4.63	4.5	8705.0	6 de abril
2	UniqOne Revlon Professional Classico Tratamien...	14.95	15.60	NaN	4.5	70.0	None
3	Revlon Professional ProYou Textura de Peinado ...	12.99	NaN	NaN	4.0	71.0	6 de abril
4	REVLON PROFESSIONAL Nutri Color Filters #400 T...	11.71	15.00	NaN	4.6	80.0	6 de abril

```
# Checking there are as many rows as products in the amazon section  
df2.shape
```

```
(796, 7)
```



Para el objetivo 2, necesitamos obtener más de 5000 artículos distribuidos en 117 páginas

1. Creamos un proyecto de scrapy y una spider.
2. Con Scrapy shell comprobamos que los XPath de los campos que nos interesan funcionan correctamente.
3. Definimos la función parse de la spider para que haga yield de un diccionario para cada producto, con todos los campos necesarios. Ponemos None por defecto si el campo está vacío para evitar errores. Dentro de la misma función, incluimos la request a la siguiente página y un callback a la función parse para que empiece de nuevo.
4. Agregamos un User Agent en Settings.
5. Hacemos crawl y lo guardamos todo en un fichero json.
6. Comprobamos que está correcto con pandas y lo convertimos a csv.

Resultado:

```
products.shape
```

```
(5601, 7)
```

```
products.head()
```

	prod_name	price_whole	price_bulk	price_old	stars	reviews	delivery_date
0	Válquer Profesional Pack Tratamiento Capilar A...	24,14	(24,14 €/unidad)	None	3,2 de 5 estrellas	27.00	mañana, 7 de abril
1	Válquer Profesional Ice Hair Mask. Mascarilla ...	13,99	(4,66 €/100 ml)	None	4,4 de 5 estrellas	24.00	mañana, 7 de abril
2	KOKEN - Champú Neutro Profesional 5L- Fórmula ...	14,90	(0,30 €/100 ml)	None	None	NaN	None
3	Pack Champu 0% Sulfatos 1000ml + Mascarilla 0%...	22,48	(22,48 €/unidad)	27,00 €	4,4 de 5 estrellas	10.00	None
4	L'Oréal Professionnel Champú Absolut Repair, 5...	11,86	(2,37 €/100 ml)	None	4,4 de 5 estrellas	4.79	domingo, 11 de abril

## FOLDERS

- ▼ amazon ○
- ▶ \_\_pycache\_\_ ○
- ▼ spiders ○
- ▶ \_\_pycache\_\_ ○
- /\* \_\_init\_\_.py ○
- /\* productos.py ○
- /\* \_\_init\_\_.py ○
- /\* items.py ○
- /\* middleware ○
- /\* pipelines.py ○
- /\* settings.py ○

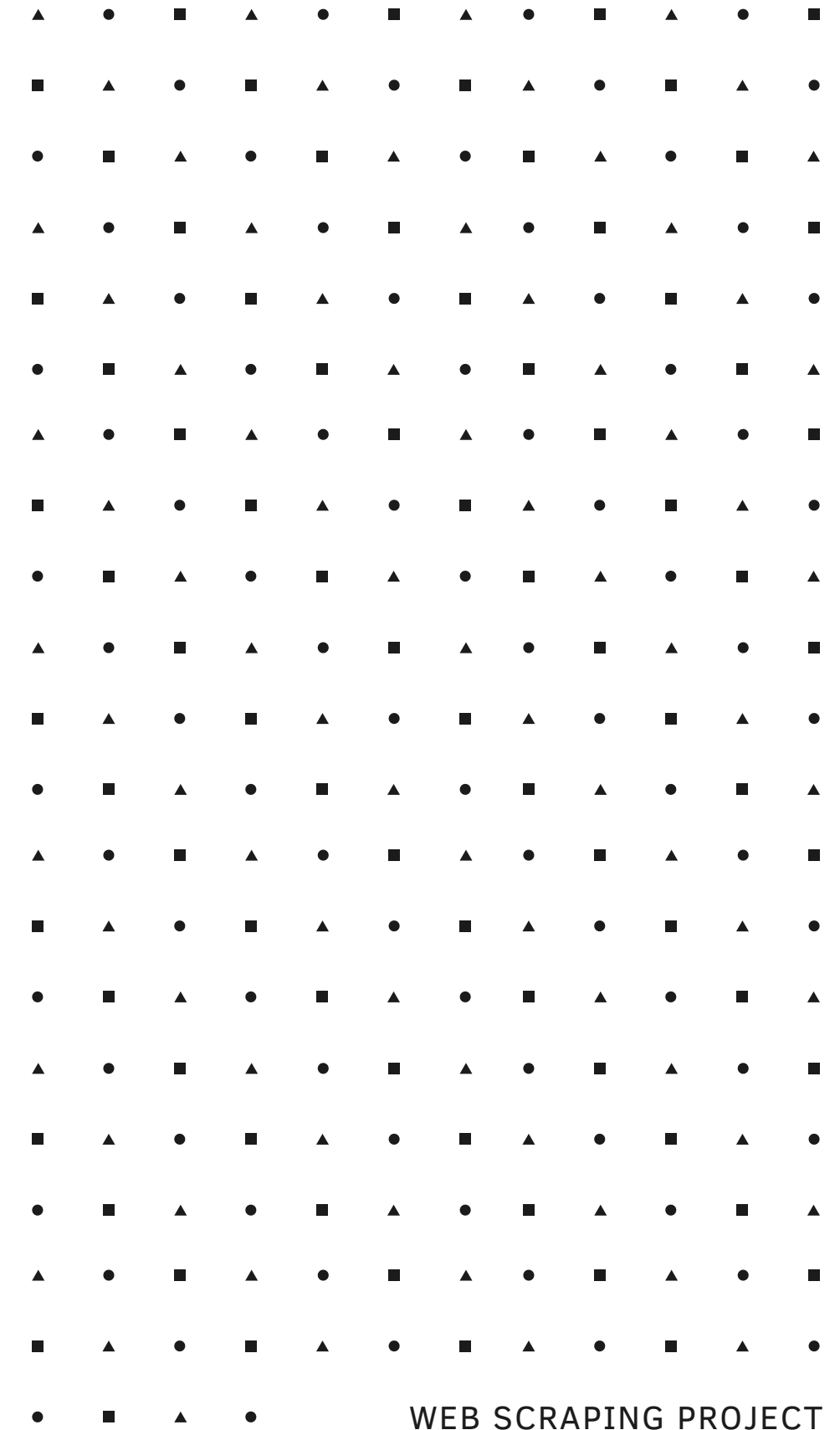
productos.py

```
1 import scrapy
2
3
4 class ProductosSpider(scrapy.Spider):
5     name = 'productos'
6     allowed_domains = ['amazon.es']
7     # Url of the HairCare Products section of Amazon Spain, with "professional" typed into the search box
8     start_urls = ['https://www.amazon.es/s?k=profesional&rh=n%3A4347699031&__mk_es_ES=%C3%85M%C3%85%C5%BD%C3%95%C3%
9
10    def parse(self, response):
11        # Obtaining a list of products by selecting them by class with xpath
12        products = response.xpath("//div[@class='sg-col-4-of-12 s-result-item s-asin sg-col-4-of-16 sg-col sg-col-4-
13        # Iterating over the list to create a dictionary for every product
14        for product in products:
15            yield {
16                'prod_name': product.xpath("../../../span[@class='a-size-base-plus a-color-base a-text-normal']/text()").get(),
17                'price_whole': product.xpath("../../../span[@class='a-price-whole']/text()").get(default=None),
18                'price_bulk': product.xpath("../../../span[@class='a-size-base a-color-secondary']/text()").get(default=None),
19                'price_old': product.xpath("../../../span[@class='a-price a-text-price']/span/text()").get(default=None),
20                'stars': product.xpath("../../../span[@class='a-icon-alt']/text()").get(default=None),
21                'reviews': product.xpath("../../../span[@class='a-size-base']/text()").get(default=None),
22                'delivery_date': product.xpath("../../../span[@class='a-text-bold']/text()").get(default=None)
23            }
24
25        # Changing to the next page when there are no more products to scrape
26        next_page = response.xpath("//li[@class='a-last']/a/@href").get()
27        if next_page is not None:
28            next_page = response.urljoin(next_page)
29            yield scrapy.Request(next_page, callback=self.parse)
```

# 04

# Obstáculos

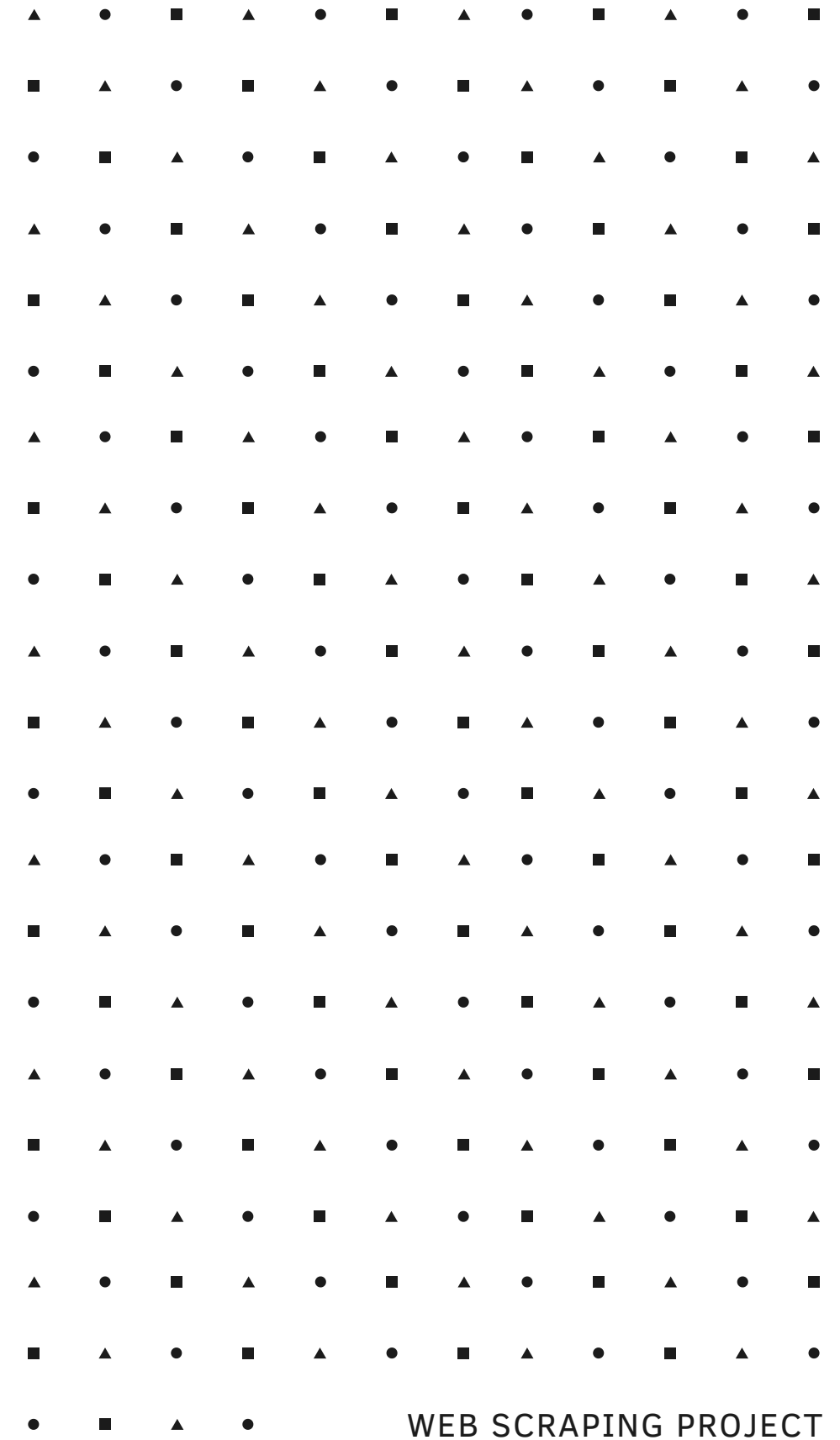
- Complejidad del código HTML de Amazon
- No existe el campo "Marca"
- En algunos productos hay campos vacíos
- Conocimiento superficial de XPath
- User Agent en Scrapy -> Settings
- Indentation error -> Convert to spaces



05

# Aprendizajes

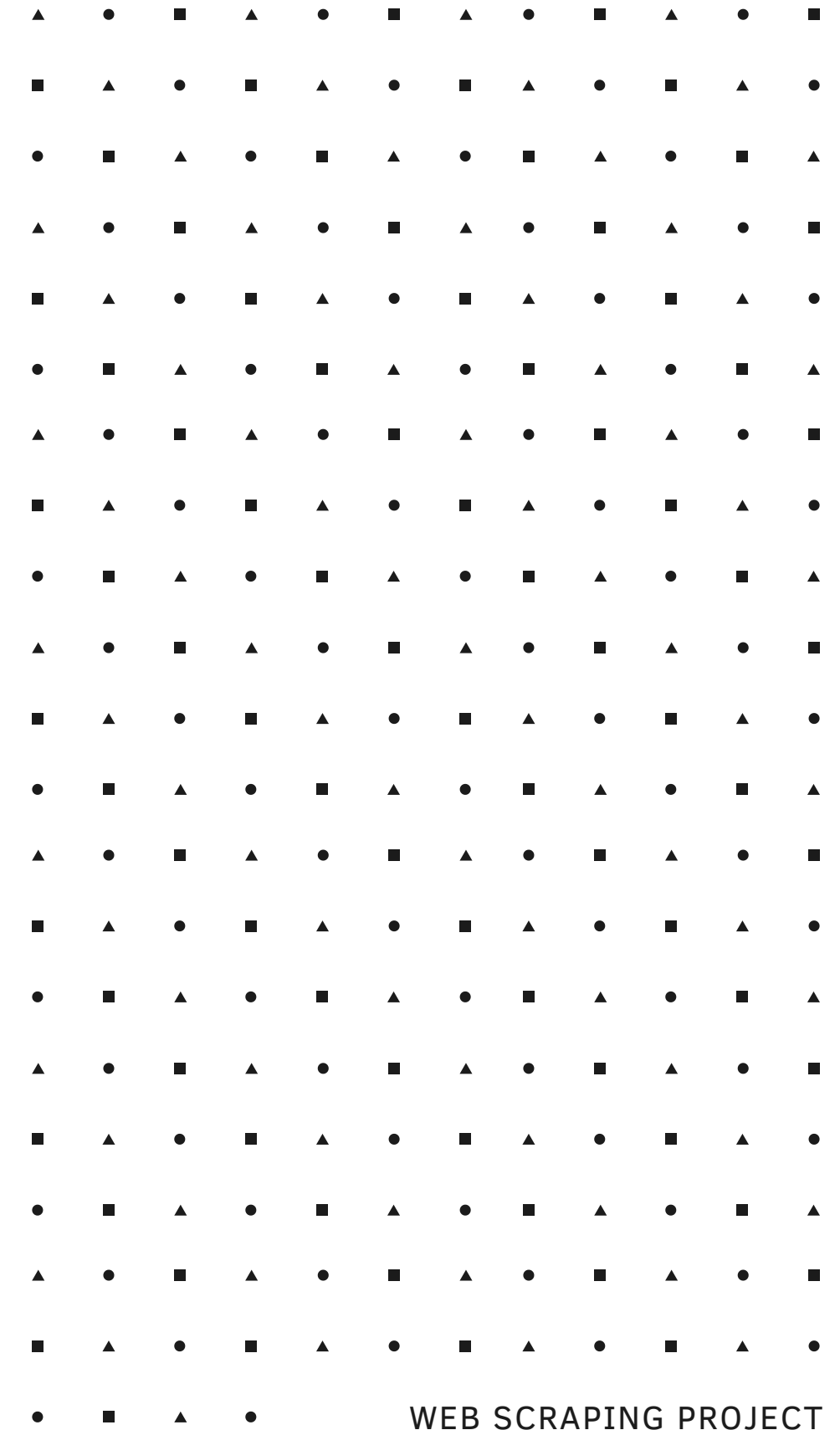
- HTML y XPath
- BeautifulSoup, Scrapy y Selenium
- Sublime Text, resolver indentation errors



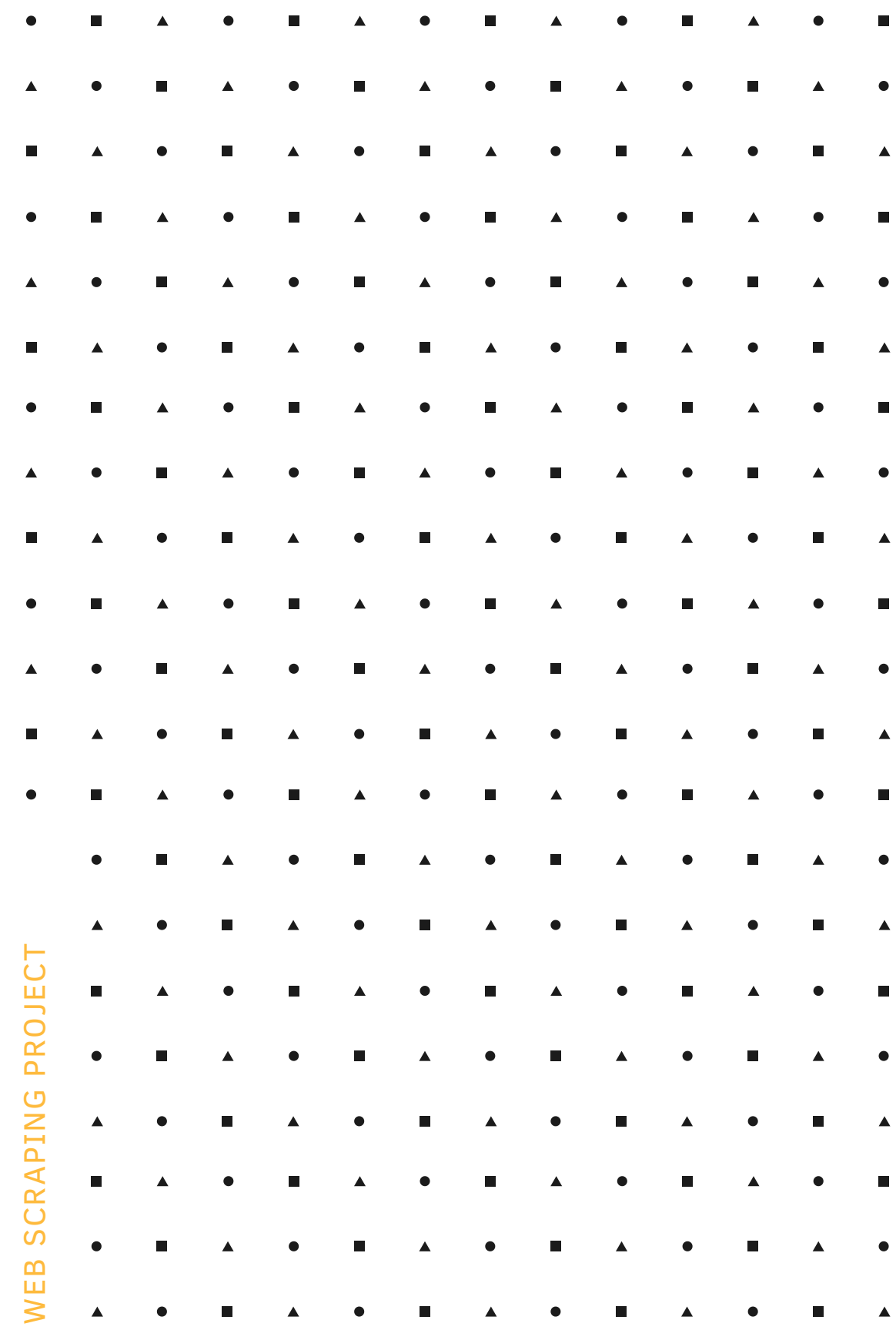
06

# Siguientes pasos

- Limpiar y analizar los datos.
- Hacer una spider más compleja que entre en la página de cada producto.
- Con más conocimiento de XPath, encontrar la marca dentro de la página del producto a partir del texto "Marca" o "Fabricante", que no siempre se encuentra dentro del mismo elemento HTML.



WEB SCRAPING PROJECT



Gracias!