

PRACTICA 1: Web Scrapping

1. Contexto. Explicar en qué contexto se ha recolectado la información.

Explique por qué el sitio web elegido proporciona dicha información.

La base de datos NANDO es publicada en una página web por la Comisión Europea y presenta la información de los organismos notificados (NB) y organismos de evaluación técnica (TAB) de acuerdo con las especificaciones técnicas armonizadas (normas armonizadas hEN y documentos de evaluación técnica EAD).

Los fabricantes europeos de materiales de construcción utilizan esta página para buscar los NB disponibles en Europa que pueden realizar la validación por tercera parte de su fabricación, como requiere el reglamento de productos de construcción RPC EU 305/2011. Esta página web se actualiza frecuentemente con los nuevos organismos notificados y aquellos que dejan de estarlo de acuerdo con los requerimientos de la legislación.

La información también resulta útil para los organismos notificados, la vigilancia de mercado y las empresas constructoras que pueden, consultando estos listados, comprobar si los organismos notificados efectivamente están autorizados para realizar esta tarea para los productos que corresponda.

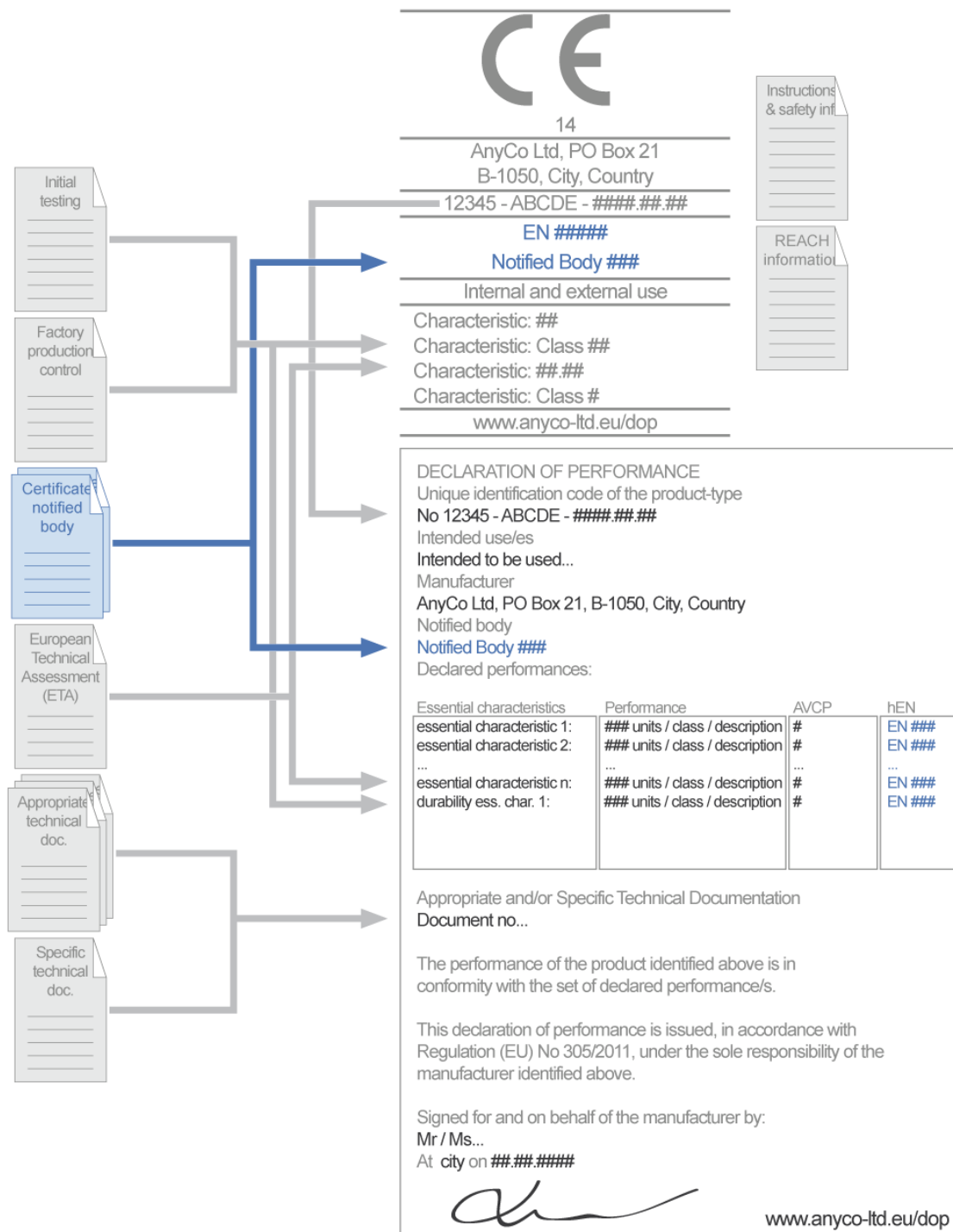
2. Definir un título para el dataset. Elegir un título que sea descriptivo.

Organismos notificados, especificaciones horizontales y normas armonizadas para la evaluación por terceras partes de materiales de construcción de acuerdo con el RPC - EU 305/2011

3. 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).

Los datos extraídos nos proporcionan una lista de las especificaciones horizontales y las normas armonizadas incluyendo sus títulos, así como de los organismos notificados para cada una de ellas y que por tanto están autorizados a realizar las tareas de validación por tercera parte que requiere el reglamento de productos de construcción RPC EU 305/2011. El listado además proporciona datos detallados de los organismos notificados

4. Representación gráfica. Presentar una imagen o esquema que identifique el dataset visualmente



5. Contenido. Explicar los campos que incluye el dataset, el periodo de tiempo de los datos y cómo se ha recogido.

- Campos:
 - Código, corresponde el código de la especificación horizontal o de la norma armonizada
 - Descripción, corresponde el título de la especificación horizontal o de la norma armonizada. Hay que destacar que la mayoría de las especificaciones horizontales no incluyen una descripción y simplemente repiten el código
 - Empresa, nombre del NB.
 - NB, código del NB.
 - País, país donde está situado el NB.
 - Teléfono, teléfono de contacto del NB.
 - Email, email de contacto del NB.
 - Web, web del NB.

- Periodo de tiempo de los datos

La lista de los NB es una foto de las organizaciones que actualmente están autorizadas a realizar la verificación de los productos de construcción. Por lo tanto, no hay un periodo exacto de validez de los datos, ya que la lista es dinámica y sufre modificaciones cuando una organización adquiere o pierde la autorización para realizar la validación.

Este script será lanzado semanalmente para asegurarnos de tener la última versión de la lista.

- Cómo se ha recogido

En la url inicial:

<https://ec.europa.eu/growth/tools-databases/nando/index.cfm?fuseaction=cp.hs&cpr=Y#hzs>

hemos iniciado un bucle que recorre la lista de las especificaciones horizontales y las normas armonizadas.

De aquí obtenemos los datos relativos a la norma armonizada o la especificación horizontal, como son el código y la descripción.

Para cada uno de los elementos de esa lista, hemos entrado en el enlace que vincula la especificación horizontal o la norma. Por ejemplo, en el primer caso TS 1187:

https://ec.europa.eu/growth/tools-databases/nando/index.cfm?fuseaction=cp.nb_hzs&hs_id=152097

En esta segunda url, hemos generado otro bucle que recorre la lista de los NB autorizados para dicha norma armonizada o especificación horizontal. Entrando en cada uno de los enlaces de cada NB. Por ejemplo, continuando con el primer caso TS 1187 y el primer NB que aparece NB 0370:

https://ec.europa.eu/growth/tools-databases/nando/index.cfm?fuseaction=cp.nb&refe_cd=EPOS_43692

Diez Alonso, Stela

Ruiz España, Silvia

Donde obtenemos los datos relativos a ese NB para añadir a nuestra base de datos (empresa, NB, país...)

Una vez recogidos todos los datos, los introducimos en el CSV.

6. Agradecimientos. Presentar al propietario del conjunto de datos. Es necesario incluir citas de investigación o análisis anteriores (si los hay).

La página web es

<https://ec.europa.eu/growth/tools-databases/nando/index.cfm?fuseaction=help.main>

El mercado CE para productos de construcción está detallado en esta guía publicada por la Comisión Europea:

https://ec.europa.eu/growth/content/ce-marking-construction-products-step-step-guide-now-available-all-eu-languages-0_en

El propietario de esta información es la Comisión Europea, el desarrollo se implementó en 2006 y esta es la guía de referencia para el uso de la página web que fue consultada durante el desarrollo:

<https://ec.europa.eu/growth/tools-databases/nando/index.cfm?fuseaction=help.main>

La guía es aplicable a todas las regulaciones y directivas europeas que requieren la designación de organismos notificados, pero en nuestro estudio sólo hemos extraído los datos correspondientes al reglamento de productos de construcción. El mismo proceso de extracción sería extrapolable a otras legislaciones con ciertas modificaciones porque hay implementaciones específicas para cada ley.

7. Inspiración. Explique por qué es interesante este conjunto de datos y qué preguntas se pretenden responder.

El problema para los usuarios de esta página web es la limitada usabilidad de la información, por ejemplo, el filtrado de organismos notificados por país debe hacerse ordenando la lista por ese campo y cuando se accede a uno de los organismos notificados deben consultarse uno a uno. Las autoridades de vigilancia de mercado deben cruzar los datos del fabricante con el correspondiente organismo notificado y especificación horizontal o norma armonizada pero este proceso no se puede realizar automáticamente sin extraer los datos previamente.

Un sistema de consulta más flexible permitiría la automatización de procesos de verificación por las organizaciones interesadas.

Esta extracción nos permitirá detectar, por ejemplo:

- Si hay suficientes NB, o si hay escasez de ellos, en cada país para una norma determinada.
- Como se comportan los NB, es decir,
 - ¿Certifican todos más o menos el mismo número de normas o hay NB “pequeños” y “grandes”?
 - ¿En todos los países presentan ese mismo comportamiento?
- Cuantos certificadores tienen las normas armonizadas y las especificaciones horizontales
 - ¿Es igual en todos los países?

- ¿Es una distribución uniforme o existen normas que pueden ser verificadas por muchos?
- ¿Hay normas que no pueden ser verificadas en un país concreto?

8. Licencia. Seleccione una de estas licencias para su dataset y explique el motivo de su selección:

El set de datos está bajo copyright de la Unión Europea, bajo licencia: CC BY-SA 4.0 License. Podemos verificar esta licencia en la siguiente página web:

https://ec.europa.eu/info/legal-notice_es

9. Código. Adjuntar el código con el que se ha generado el dataset, preferiblemente en Python o, alternativamente, en R.

Nota: El código está incluido en el repositorio Github.

```

1  # Importamos las librerías que vamos a necesitar
2  import re
3  import requests
4  from bs4 import BeautifulSoup
5  import csv
6  import time
7  import unicodedata
8
9  # Función para eliminar tildes
10 def eliminar_tildes(cadena):
11     s = ''.join((c for c in unicodedata.normalize('NFD',cadena) if unicodedata.category(c) != 'Mn'))
12     return s
13
14 # Creamos el documento .csv para guardar los campos de interés
15 with open("Standards.csv","w", encoding="utf-8") as csvfile:
16     spamwriter=csv.writer(csvfile, delimiter=',')
17     # Cabecera
18     spamwriter.writerow(['CODIGO','DESCRIPCION', 'EMPRESA', 'CODIGO_NB', 'PAIS_NB', 'TELEFONO_NB', 'EMAIL_NB', 'WEB'])
19     cnt=0
20     Datos=[]
21
22     # Obtenemos los datos
23     #Calculamos el tiempo de respuesta para introducir retrasos en las peticiones consecutivas
24     for term in ["web scrapping", "web crawling", "scrape this site"]:
25         t0 = time.time()
26         r = requests.get("https://ec.europa.eu/growth/tools-databases/nando/index.cfm?fuseaction=cp.hs&cpr=Y", para
27
28     ## Calculo final de la estimación del tiempo de respuesta en segundos
29     response_delay = time.time() - t0
30
31
32 # Marcamos la url base para los enlaces siguientes extraidos de la web
33 url_base = 'https://ec.europa.eu/growth/tools-databases/nando/'
34
35 # Listas para guardar los campos de interés
36 Codigo=[]
37 Descripcion=[]
38 Empresa=[]
39 nb=[]
40 Country=[]
41 Phone=[]
42 Email=[]
43 Web=[]
44
45 # Cargamos la primera página, donde se encuentran el listado de normas armonizadas.
46 page = requests.get('https://ec.europa.eu/growth/tools-databases/nando/index.cfm?fuseaction=cp.hs&cpr=Y')
47 soup = BeautifulSoup(page.content)
48
49 contador=-1
50
51 # Recorremos la lista de normas:
52 for i in soup.find_all('a',{'class':'list'}):
53     # Cada norma introducimos un tiempo de espera para espaciar nuestras peticiones
54     # Espera de 10 veces el tiempo de respuesta
55     time.sleep(10 * response_delay)
56
57     #####
58     ### CAMPO CODIGO ###
59     #####
60     codigo = i.string # Valor a introducir en el campo: Codigo
61
62     #####
63     ### CAMPO DESCRIPCION ###
64     #####

```

```

64
65     cod_simple = re.findall('[A-Z]+ [0-9]*[-]*[0-9]*[:]*[0-9]*[/]*[*]*[A-Z]*[:]*[0-9]*', codigo)
66     if cod_simple[0] == 'EN ISO':
67         cod_simple = re.findall('[A-Z]+ [A-Z]* [0-9]*[-]*[0-9]*[:]*[0-9]*[/]*[*]*[A-Z]*[:]*[0-9]*', codigo)
68     descripcion = soup.find_all(string=re.compile('\t'+cod_simple[0]))
69
70     try:
71         descripcion = (descripcion[0].strip()) # Valor a introducir en el campo: Descripcion
72     except IndexError:
73         descripcion = "---"
74     # print(descripcion)
75
76     if contador==1:
77         break
78     else:
79         contador=contador+1
80
81     # Cargamos la segunda página, de esta página no extraeremos datos con los que rellenar campos
82     # pero muestra un listado de los NB que pueden realizar la notificación de la norma que estamos
83     # tratando.
84     page2 = requests.get(url_base+i.get('href'))
85     soup2 = BeautifulSoup(page2.content)
86
87     # Recorremos el listado de NB de la norma
88     for j in soup2.find_all('a',{'class':'list'}):
89         #Para cada empresa introducimos un tiempo de espera para espaciar nuestras peticiones
90         #Espera de 2 veces el tiempo de respuesta
91         time.sleep(2 * response_delay)
92
93         #Para cada elemento de la segunda página, es decir, para cada NB cargamos una página
94         #en la que están sus datos.
95         page3 = requests.get(url_base+j.get('href'))
96         soup3 = BeautifulSoup(page3.content)
97
98         #####
99         ### CAMPO EMPRESA y WEB ###
100         #####
101         for k in soup3.find_all('strong'):
102             empresa=(k.string) # Valor a introducir en el campo: Empresa
103             empresa=elimina_tildes(empresa)
104             Empresa.append(empresa)
105
106             web=soup3.find('a',{'target':'_blank'}) # Valor a introducir en el campo: Web_NB
107             if web:
108                 # print(web.get('href'))
109                 Web.append(web.get('href'))
110             else:
111                 Web.append("---")
112
113         #####
114         ### CAMPOS NB, PAIS, TELEFONO, EMAIL ###
115         #####
116         NB = soup3.find_all(string=re.compile("Notified Body number :"))
117         aux = NB[0].replace('Notified Body number :', '')
118         NB = aux.strip() #Valor a introducir en el campo:Codigo_NB
119         nb.append(NB)
120
121         country = soup3.find_all(string=re.compile("Country :"))
122         aux = country[0].replace('Country :', '')
123         country = aux.strip() #Valor a introducir en el campo: Pais_NB
124         Country.append(country)
125
126         phone = soup3.find_all(string=re.compile("Phone :"))
127         aux = phone[0].replace('Phone :', '')
128         phone = aux.strip() #Valor a introducir en el campo: Telefono_NB
129         Phone.append(phone)
130
131         email = soup3.find_all(string=re.compile("Email :"))
132         aux = email[0].replace('Email :', '')
133         email = aux.strip() #Valor a introducir en el campo: Email_NB
134         Email.append(email)
135
136         # Almacenamos el Código y la Descripción
137        Codigo.append(codigo)
138         Descripcion.append(descripcion)
139
140         #####
141         ### GUARDAMOS LOS DATOS ###
142         #####
143         # Montamos los datos en una fila y lo guardamos en el csv
144         Rows=[Codigo[cnt], Descripcion[cnt], Empresa[cnt], nb[cnt], Country[cnt], Phone[cnt], Email[cnt], Web[cnt]]
145         Datos.append(Rows)
146         # print(Datos[cnt])
147         spamwriter.writerow(Datos[cnt])
148
149         #Imprimimos un aviso cada 25 registros para saber que el proceso continua.
150         if cnt%25 == 0:
151             print("Guardando fila... " + str(cnt))
152             cnt=cnt+1
153
154     csvfile.close()
155     print("Proceso finalizado. \nSe han guardado "+str(cnt)+" registros")
156

```

10. Dataset. Publicación del dataset en formato CSV en Zenodo con una pequeña descripción.

Título: Notified bodies, horizontal specifications and harmonised standards for construction materials third party assessment according to the CPR - EU 305/2011

Descripción: This dataset was extracted from NANDO, the website of the European Commission for notified bodies. It contains the list of notified bodies, horizontal specifications and harmonised standards for construction materials third party assessment according to the CPR - EU 305/2011

Enlace: <https://zenodo.org/record/3741434#.XosBNnLtZEY>

DOI: <https://doi.org/10.5281/zenodo.3741434>

11. Entrega. Presentar el trabajo con el DOI del dataset en Github.

DOI del dataset de Github: <https://zenodo.org/badge/latestdoi/253556749>

Enlace: <https://github.com/Stela-Silvia/Standards-scraping/tree/1.0>

Documento realizado por:

Diez Alonso, Stela

Ruiz España, Silvia

Contribuciones	Firma
Investigación previa	S.D.A y S.R.E.
Redacción de las respuestas	S.D.A y S.R.E.
Desarrollo código	S.D.A y S.R.E.