

SPYDER

Rafael Fernández Ortiz

Se denomina **Spyder** o **Web Spyder** al bot/script que se utiliza para analizar o rasrear el código html de una dirección url. De este modo, puedes identificar las palabras y por tanto hacer un análisis de frecuencia de ciertas palabras, o identificar los distintos hipervínculos que existen en dicha url.

Este es un primer proyecto de clase, así repaso el concepto y la utilidad de las clases, de un objeto *spyder*.

Las distintas **funciones** o **métodos** del *spyder* permiten, o bien extraer los links de direcciones *url* y *pdfs*, y mostrarlo en pantalla mediante una lista; o bien, escribirlo en un archivo *.txt*

Comentarios: Este es un primer esbozo del código, aun hay que afinarlo muchísimo, pero como primera instance proporciona buenos resultados.

```
import urllib.request as url
class Spyder(object):

    def __init__(self,url):
        self.url = str(url)

    # Funciones privadas previas necesarias para las funciones publicas
    def __getHtml(self):
        dir = url.urlopen(self.url)
        html = str(dir.read())
        return html
    def __getUrIs(self):
        urIs = []
        html = self.__getHtml()
        while True:
            pos = html.find('href="')
            if pos == -1:
                break
            html = html[(pos + 6):]
            urIs.append(html[:html.find('"')])
        return urIs
    def __busca(self,cosa):
        return [x for x in self.__getUrIs() if str(cosa) in x]

    # Funciones publicas
    def getLink(self): #Devuelve una lista con los enlaces en la url
        return self.__busca("http")
    def getPdf(self): #Devuelve una lista con los enlaces de los pdfs adjuntos
        return self.__busca(".pdf")
    def getLinktoTxt(self): #Crea un .txt con los enlaces que hay en la url
        fich = open("links.txt", 'w')
        link = self.getLink()
        for i in link:
            fich.write(str(i)+"\n")
```

```
fich.close()
return fich
def getPdfToTxt(self): #Crea un .txt con los enlaces de los pdfs adjuntos
    fich = open("pdfs.txt", 'w')
    link = self.getPdf()
    for i in link:
        fich.write(str(i)+"\n")
    fich.close()
    return fich
```

##-----

#Ejemplo

```
sp1 = Spyder("https://es.wikipedia.org/wiki/Madrid")
sp1.getPdfToTxt()
sp1.getLinkToTxt()
```