

Ing. Sistemas Computaciones

"Obtención de Información semi-pasiva"

Materia

GOBIERNO DE TI|COMPUTER SECURITY

Grupo 6CV2

Maestro
Aldama Coahuila Mario Alberto

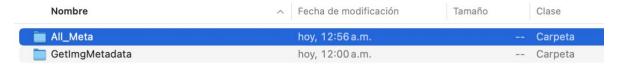
Autor: Aldo Alcántara Martínez

Boleta 2019630578

Fecha: 22/03/2024

Reporte de los pasos que realice:

1.- Lo primero es designar una carpeta para trabajar nuestro script y no mezclarlo, en mi caso yo la llame All Meta



2.- Despues, abro mi terminal y me posiciono en la carpeta que acabo de crear, luego procedo a crear una carpeta venv y así como se ve en la siguiente linea activo el modo venv, que ya en la tercer linea se ve como se activa en la parte izquierda, por ultimo reviso mi php List y compruebo que es lo que tengo instalado.

source venv/bin/activate: Esta instrucción activa el entorno virtual creado en el paso anterior. Una vez activado, cualquier instalación de paquetes de Python se realizará en el entorno virtual y no afectará al sistema global de Python.

3.- Ya que comprobémoste que tengo instalado, ahora instalo las bibliotecas que voy a ocupar, que son las que muestro en la imagen, junto con su proceso de descarga y por ultimo compruebo que se instalaron.

pip install python-docx openpyxl PyPDF2: Esta
instrucción instala los paquetes necesarios para el
script. python-docx se utiliza para trabajar con
documentos DOCX, openpyxl para trabajar con archivos
XLSX y PyPDF2 para trabajar con archivos PDF

```
(venv) aldoalcantara@MacBook-Pro-de-Aldo-2 All_Meta % pip install python-docx openpyxl PyPDF2
Collecting python-docx
  Downloading python_docx-1.1.0-py3-none-any.whl.metadata (2.0 kB)
Collecting openpyxl
  Downloading openpyxl-3.1.2-py2.py3-none-any.whl.metadata (2.5 kB)
Collecting PyPDF2
Downloading pypdf2-3.0.1-py3-none-any.whl.metadata (6.8 kB) Collecting lxml>=3.1.0 (from python-docx)
  Downloading lxml-5.1.0-cp312-cp312-macosx_10_9_x86_64.whl.metadata (3.5 kB)
Collecting typing-extensions (from python-docx)
Downloading typing_extensions-4.10.0-py3-none-any.whl.metadata (3.0 kB)
Collecting et-xmlfile (from openpyxl)
  Downloading et_xmlfile-1.1.0-py3-none-any.whl.metadata (1.8 kB)
Downloading python_docx-1.1.0-py3-none-any.whl (239 kB)
                                                239.6/239.6 kB 2.1 MB/s eta 0:00:00
Downloading openpyxl-3.1.2-py2.py3-none-any.whl (249 kB)
                                                 250.0/250.0 kB 6.1 MB/s eta 0:00:00
Downloading pypdf2-3.0.1-py3-none-any.whl (232 kB)
                                                232.6/232.6 kB 7.4 MB/s eta 0:00:00
Downloading lxml-5.1.0-cp312-cp312-macosx_10_9_x86_64.whl (4.8 MB)
                                                 4.8/4.8 MB 24.5 MB/s eta 0:00:00
Downloading et_xmlfile-1.1.0-py3-none-any.whl (4.7 kB)
Downloading typing_extensions-4.10.0-py3-none-any.whl (33 kB)
Installing collected packages: typing-extensions, PyDPF2, lxml, et-xmlfile, python-docx, openpyxl Successfully installed PyPDF2-3.0.1 et-xmlfile-1.1.0 lxml-5.1.0 openpyxl-3.1.2 python-docx-1.1.0 typing-extensions-4.10.0
(venv) aldoalcantara@MacBook-Pro-de-Aldo-2 All_Meta % pip list
Package
                   Version
et-xmlfile
lxml
                   5.1.0
openpyxl
                   3.1.2
                    24.0
pip
PyPDF2
                    3.0.1
python-docx
                   1.1.0
typing_extensions 4.10.0
4.- Por ultimo procedo a correr el programa y verifico
```

que todo funciono correctamente

```
(venv) aldoalcantara@MacBook-Pro-de-Aldo-2 All_Meta % python all_metadata_info.py
Metadatos del archivo SegundoMilestone.docx:
 author: Jhoana Monserrat Pimentel Lopez
  title:
 created: 2023-11-27 16:33:00
 modified: 2023-11-27 16:33:00
Metadatos del archivo Horario.xlsx:
 author: None
  title: None
 created: 2024-01-30 19:05:09
 modified: 2024-03-22 06:00:13
Metadatos del archivo Horario_Oficial.pdf:
  author: None
  creator: Crystal Reports
 producer: Powered By Crystal
 title: None
 subject: None
 created: None
 modified: None
```

Codigo

```
import os
import re
from docx import Document
from openpyxl import load workbook
from PyPDF2 import PdfReader
def extract metadata from pdf(filename):
    pdf file = PdfReader(filename)
   metadata = {
        "author": pdf file.metadata.author,
        "creator": pdf file.metadata.creator,
        "producer": pdf file.metadata.producer,
        "title": pdf file.metadata.title,
        "subject": pdf file.metadata.subject,
        "created": pdf_file.metadata.get('/CreationDate', None),
        "modified": pdf file.metadata.get('/ModDate', None)
    return metadata
def extract metadata from docx(filename):
    document = Document(filename)
    core properties = document.core properties
    metadata = {
        "author": core properties.author,
        "title": core properties.title,
        "created": core properties.created,
        "modified": core_properties.modified
    return metadata
def extract metadata from xlsx(filename):
    wb = load workbook(filename)
   author = None
    if wb.sheetnames: # Verifica si hay hojas de cálculo en el libro
        sheet = wb[wb.sheetnames[0]] # Obtiene la primera hoja de
cálculo
```

```
if sheet.dimensions: # Verifica si la hoja de cálculo contiene
datos
            author = sheet.cell(row=1, column=1).value # Obtiene el
valor de la primera celda
    metadata = {
        "author": author,
        "title": wb.properties.title,
        "created": wb.properties.created,
        "modified": wb.properties.modified
    return metadata
def process files(directory path):
    for file name in os.listdir(directory path):
        file path = os.path.join(directory path, file name)
        if os.path.isfile(file path):
            file extension = re.findall(r"\.(pdf|docx|xlsx)$", file name)
            if file extension:
                file extension = file extension[0]
                if file extension == "pdf":
                    metadata = extract_metadata_from_pdf(file_path)
                elif file extension == "docx":
                    metadata = extract metadata from docx(file path)
                elif file extension == "xlsx":
                    metadata = extract metadata from xlsx(file path)
                print(f"Metadatos del archivo {file name}:")
                for k, v in metadata.items():
                    print(f" {k}: {v}")
                print()
# Ruta al directorio
directory_path = "/Users/aldoalcantara/Documents/ESCOM/Computer
Security/Python/All Meta"
# Procesar los archivos en el directorio
process files(directory path)
```