

ESCUELA POLITÉCNICA NACIONAL

ESCUELA DE FORMACIÓN DE TECNÓLOGOS

**DESARROLLO DE APLICACIÓN WEB PARA EXTRACCIÓN DE
DATOS CON WEB SCRAPING Y HERRAMIENTA DE ANALÍTICA
DE DATOS**

APLICATIVO WEB DE EXTRACCIÓN DE DATOS.

**TRABAJO DE INTEGRACIÓN CURRICULAR PRESENTADO COMO
REQUISITO PARA LA OBTENCIÓN DEL TÍTULO DE TECNÓLOGO SUPERIOR
EN DESARROLLO DE SOFTWARE**

CHRISTIAN GABRIEL SOLEDISPA MENDOZA

DIRECTOR: JUAN PABLO ZALDUMBIDE PROAÑO

DMQ, Agosto 2023

CERTIFICACIONES

Yo, Christian Gabriel Soledispa Mendoza declaro que el trabajo de integración curricular aquí descrito es de mi autoría; que no ha sido previamente presentado para ningún grado o calificación profesional; y, que he consultado las referencias bibliográficas que se incluyen en este documento.



CHRISTIAN GABRIEL SOLEDISPA MENDOZA

christian.soledispa@epn.edu.ec

cgsoledispa@outlook.com

Certifico que el presente trabajo de integración curricular fue desarrollado por Christian Gabriel Soledispa Mendoza, bajo mi supervisión.



ING. JUAN PABLO ZALDUMBIDE PROAÑO

Juan.zaldumbide@epn.edu.ec

DECLARACIÓN DE AUTORÍA

A través de la presente declaración, afirmamos que el trabajo de integración curricular aquí descrito, así como el producto resultante del mismo, son públicos y estarán a disposición de la comunidad a través del repositorio institucional de la Escuela Politécnica Nacional; sin embargo, la titularidad de los derechos patrimoniales nos corresponde a los autores que hemos contribuido en el desarrollo del presente trabajo; observando para el efecto las disposiciones establecidas por el órgano competente en propiedad intelectual, la normativa interna y demás normas.

CHRISTIAN GABRIEL SOLEDISPA MENDOZA

DEDICATORIA

El presente proyecto va dedicado a toda la comunidad de estudiantes de la Escuela Politécnica Nacional conformada por excelentes docentes y estudiantes que cada día se esfuerzan por alcanzar el éxito y sobre todo por el progreso de la sociedad. También a las personas a las que el siguiente proyecto para la extracción de datos con webscraping les sea de utilidad.

SOLEDISPA MENDOZA CHRISTIAN GABRIEL

AGRADECIMIENTO

Agradezco a mis padres y hermanos por ser mi principal apoyo durante todo este proceso en mi vida estudiantil, porque sin ellos no podría alcanzar todo lo que he logrado hasta ahora con su amor incondicional y valores que me supieron enseñar este triunfo más que para mí es para ellos.

A mi amiga Viviana por haber estado incondicionalmente a mi lado desde el primer semestre que estuve en la EPN, acompañarme en este largo y sacrificado proceso con su amistad.

A mi amiga Doménica por ser mi apoyo moral e incondicional este último tiempo que ha sido muy difícil, pero con su amistad y alegría me ha llenado de fuerzas para seguir adelante y alcanzar mis objetivos.

SOLEDISPA MENDOZA CHRISTIAN GABRIEL

ÍNDICE DE CONTENIDO

CERTIFICACIONES.....	I
DECLARACIÓN DE AUTORÍA	II
DEDICATORIA.....	III
AGRADECIMIENTO.....	IV
ÍNDICE DE CONTENIDO.....	V
RESUMEN	VII
ABSTRACT	VIII
1 DESCRIPCIÓN DEL COMPONENTE DESARROLLADO.....	1
1.1 Objetivo general.....	2
1.2 Objetivos específicos	2
1.3 Alcance	2
1.4 Marco Teórico	3
Metodología	3
Web scraping.....	3
Sistema web	3
2 METODOLOGÍA	4
2.1 Metodología de Desarrollo	4
2.2 Diseño de la arquitectura	5
Arquitectura de Datos	5
Patrón arquitectónico Modelo Vista Controlador	6
2.3 Herramientas de Desarrollo	7
3 RESULTADOS.....	9
<i>Sprint 0.</i> Configuración del ambiente de desarrollo.....	9
<i>Sprint 1.</i> Elección del lenguaje de programación.....	10
<i>Sprint 2.</i> Cambio del lenguaje de programación	12
<i>Sprint 3.</i> Codificación de Scripts en pyhton para web scraping.....	13
<i>Sprint 4.</i> Implementación de registro de usuarios.....	20
<i>Sprint 5.</i> Implementación de Lista de usuarios.....	21
<i>Sprint 6.</i> Implementación de Actualización de usuarios.....	22
<i>Sprint 7.</i> Implementación de Borrar usuarios.....	23
<i>Sprint 8.</i> Implementación Login para el usuario en la web.....	24
<i>Sprint 9.</i> Implementación de scraping en la interfaz web	25

<i>Sprint 10.</i> Almacenamiento de usuarios registrados en la base de datos	26
<i>Sprint 11.</i> Scraping anchors en la web.....	27
<i>Sprint 12.</i> Rutas de Django y Modelos de usuario.....	29
<i>Sprint 13.</i> Implementación de vistas	31
<i>Sprint 14.</i> Implementación de API Rest de autenticación y API Rest de listas de usuarios	32
<i>Sprint 15.</i> Implementación total del web scraping.....	34
<i>Sprint 16.</i> Almacenamiento de los datos del historial en la base de datos	36
<i>Sprint 17.</i> Eliminación de la API Rest y base de datos	38
4 Conclusiones	40
5 Recomendaciones	41
6 REFERENCIAS BIBLIOGRÁFICAS	42
7 ANEXOS	44
ANEXO I Certificado turnitin.....	45
ANEXO II Manual Técnico	46
ANEXO III Manual de usuario	51
ANEXO IV Manual de instalación.....	52

RESUMEN

Actualmente existen sistemas web que permiten realizar web scraping pero son difíciles de buscar o encontrar en el internet, esta funcionalidad es muy útil en muchos ámbitos tanto de estudio como laborales. Este sistema web para la extracción de datos tiene como objetivo facilitar este proceso y almacenar la información extraída de manera ordenada y estructurada en una base de datos NoSQL como lo es MongoDB.

¿Por qué es importante saber utilizarla si necesitas analizar datos de la web? De esta manera automatizada se puede obtener información de páginas que lo permitan y estructurar dicha información para que sea utilizada a conveniencia de los usuarios. Al aplicar esta herramienta la obtención de datos facilita la visualización de estos mediante un sistema de frontend.

Una herramienta muy importante para el desarrollo de este sistema web es Django ya que es enfocado al lenguaje de programación python y as vez este lenguaje es el más práctico para realizar web scraping de diferentes páginas web y estructurarlos en una base de datos que este caso es MongoDB.

PALABRAS CLAVE: web scraping, NoSQL, Django, Python, MongoDB.

ABSTRACT

Currently there are web systems that allow web scraping to be carried out but they are difficult to search or find on the internet, this functionality is very useful in many areas of both study and work. This web system for data extraction aims to facilitate this process and store the extracted information in an orderly and structured manner in a NoSQL database such as MongoDB.

Why is it important to know how to use it if you need to analyze web data? In this automated way, information can be obtained from pages that allow it and structure said information so that it can be used at the convenience of users. By applying this tool, data collection facilitates its visualization through a frontend system.

A very important tool for the development of this web system is Django since it is focused on the python programming language and at the same time this language is the most practical to carry out web scraping of different web pages and structure them in a database that in this case is MongoDB.

KEYWORDS: web scraping, NoSQL, Django, Python, MongoDB.

1 DESCRIPCIÓN DEL COMPONENTE DESARROLLADO

El sistema web permite la extracción de datos de diferentes páginas web con web scraping, los datos son almacenados en una base de datos NoSQL que en este caso es MongoDB.

El web scraping es una herramienta muy útil para recopilar todo tipo de datos que se encuentren online y al almacenarlos poder utilizarlos y manipularlos de la manera que sea conveniente y cualquier fin que tenga el usuario [1].

Se puede scrapear todo tipo de datos que se encuentren en la web, la mayoría de las páginas web disponen sus datos para que los scrapers puedan extraerlos [2]. El web scraping tiene muchos propósitos como realizar estudios de mercado, buscar potenciales clientes, extraer información valiosa como noticias y contenido que sea del interés del usuario [1].

Un buen webscraping es capaz de extraer datos de forma fiable tomando en cuenta aspectos importantes como utilizar localizadores precisos para poder especificar lo que la información que se quiere extraer exactamente. La calidad de los datos no es prudente trabajar con datos no estructurados ya que serpia desordenado y no se puede ofrecer buenos resultados en la búsqueda por eso se debe asegurar que los datos estén clasificados antes de la entrega [1].

El web scraping se relaciona con la indexación de la web, que indexa dicha información de la web mediante un robot lo que es utilizado por los motores de búsqueda. El web scraping se enfoca más en la transformación de los datos sin estructura en la web en datos estructurados que se pueden almacenar y analizar en bases de datos centrales, hojas de cálculo o en cualquier fuente de almacenamiento de datos [3] como se realiza en el siguiente proyecto de web scraping con la base de datos NoSql como MongoDB.

1.1 Objetivo general

Desarrollar un aplicativo web para la extracción de datos.

1.2 Objetivos específicos

1. OE1: Determinar los requerimientos globales.
2. OE2: Diseñar el modelo de la base de datos y la arquitectura
3. OE3: Extraer los datos de diferentes páginas web con web scraping.
4. OE4: Almacenar los datos extraídos en una base de datos de tiempo real.
5. OE5: Realizar la estructura de datos para que de esta manera puedan ser visualizados.

1.3 Alcance

El presente proyecto propone el despliegue de un sistema web para la extracción de datos con web scraping, por medio de herramientas funcionales como Java y Python para escribir la lógica del lado del servidor y tecnologías de API para permitir la comunicación entre sistemas que en la mayor parte del desarrollo se busque resultados precisos y claros.

El sistema proveerá de dos perfiles; el perfil administrador desde el cual se podrá administrar y gestionar la información extraída y realizará las conexiones con la base de datos en tiempo real y un perfil usuario específico donde los usuarios ingresan y buscan la información que desea obtener.

Este sistema permitirá al usuario acceder a información específica de ciertas páginas web como palabras más buscadas, tendencias e información precisa que se quiera buscar de ciertas publicaciones realizadas que con las herramientas que se tienen a mano es muy complicado de conseguir. El sistema de extracción se encargará únicamente de extraer, gestionar y analizar la información para almacenarlas en una base de datos en tiempo real y luego esta información pueda ser visualizada mediante el componente de frontend.

1.4 Marco Teórico

Metodología

Una metodología de desarrollo de software es un conjunto de métodos y técnicas utilizadas para diseñar y desarrollar un software informático. Por cuestión de organización es indispensable trabajar con una metodología para tener todos los factores ordenados y saber cómo utilizarlos, también sirven para que el desarrollo del proyecto tenga un control y así minimizar los márgenes de errores. Utilizar una metodología en el desarrollo hace que se gestionen de mejor manera los recursos y también ahorra tiempo [4].

Web scraping

Mediante el web scraping se extraen y almacenan todo tipo de datos de páginas web para que luego sean analizados y utilizados como sea conveniente. Todos los datos que sean extraídos se pueden almacenar en bases de datos. El web scraping es utilizado para muchas tareas como la recopilación rápida de datos, en el ámbito profesional para obtener ventajas realizando un estudio de mercado de la competencia y así comparar con los productos propios y por esta razón es posible bloquearlo como medio de seguridad mediante bots de software [5].

Sistema web

Un Sistema web es un software al que se accede por medio de un servidor web sin la necesidad de tener una aplicación de escritorio. En la actualidad el sistema web es muy utilizado ya que es muy práctico y rápido en la web. Las aplicaciones web evitan gastos para que no sea necesario el uso de programas que por lo general son muy costosos y se puede trabajar desde cualquier lugar donde se encuentre. Los sistemas web al trabajar con bases de datos pueden procesar y mostrar la información a los usuarios, estos sistemas no tienen la necesidad de ser instalados [6].

2 METODOLOGÍA

En la programación se utilizan las metodologías de desarrollo de software, el principal objetivo de estas refiriéndose a este ámbito es poder realizar un trabajo en equipo de manera organizada y con el paso del tiempo han pasado de ser únicamente para la organización a ser una base indispensable al momento de desarrollar software. Últimamente las metodologías ágiles de desarrollo software se han impuesto sobre muchas otras por lo tanto es muy importante saberlas para poder desarrollar software de calidad y eficaz [4]. El caso de estudio de este proyecto para la extracción de datos por medio de webscraping es aplicar una metodología de software ágil.

2.1 Metodología de Desarrollo

RAD (Rapid application development) es una metodología de desarrollo de software que se basa en prototipos, comentarios y pruebas rápidas en lugar de una planificación en específico. Con un desarrollo rápido de aplicaciones, los desarrolladores realizan varias iteraciones y actualizaciones del software sin la necesidad de empezar desde el inicio lo que garantiza que el principal objetivo sea la calidad del software y solvente los requisitos de los usuarios finales [5].

Al hacer que los usuarios puedan interactuar con prototipos que se encuentran en constante evolución la funcionalidad comercial del proyecto puede ser mucho mayor a la que se logra mediante el modelo cascada. Con RAD los errores que puedan ocurrir se pueden descubrir y actuar sobre estos en etapas tempranas del desarrollo [6].

RAD es la mejor manera para desarrollar prototipos de manera rápida para probar cada una de las funcionalidades del software sin tener que preocuparse por el producto final [5].

Recopilación de Requerimientos

Un requerimiento es una capacidad o una condición que un usuario necesita para lograr un objetivo o resolver algún problema, son capacidades o condiciones que debe tener un sistema para solventar las necesidades específicas que tenga el usuario. Es un proceso de documentación, definición y mantenimiento de los requisitos, son un conjunto de

actividades que se encargan de identificar y transmitir el propósito que tiene un sistema o software y el contexto en que será utilizado [10].

2.2 Diseño de la arquitectura

En el diseño arquitectónico es una planificación basada en modelos, abstracciones teóricas y patrones que se toman en cuenta al momento de realizar un software con cierta complejidad como un paso previo a cualquier implementación de este. Así se obtiene una guía teórica detallada que permite armar de manera correcta dicho software [11]. A continuación, se presenta el modelo arquitectónico del que se basa el respectivo desarrollo del proyecto.

Arquitectura de Datos

En este proyecto entra Django que es un framework de código abierto escrito en Python de alto nivel con el que se fomenta el desarrollo limpio, práctico y rápido [12]. Django es una herramienta utilizada para desarrollar tanto full-stack de páginas web y aplicaciones como el desarrollo de servidores. Es el mejor framework para desarrollar aplicaciones web con Python y es uno de los más demandados por los desarrolladores que trabajan con este lenguaje en específico [13].

En la **Fig. 1** se observa en breve descripción las herramientas que forman parte del proyecto en desarrollo, y además una breve visualización del Modelo Vista Controlador (MVC).

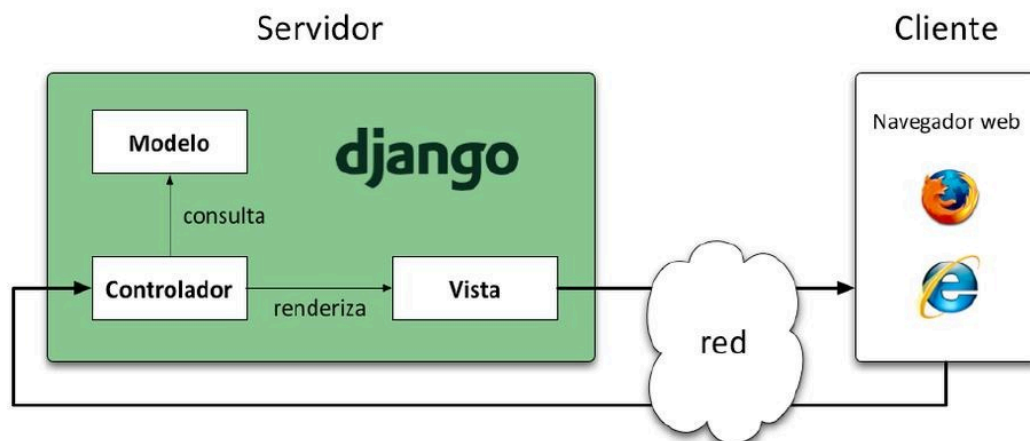


Fig. 1: Patrón arquitectónico

Patrón arquitectónico Modelo Vista Controlador

El MVC es una arquitectura de software que se utiliza para separa el código en cada una de sus funcionalidades en distintas capas las cuales se encargan de realizar una funcionalidad en específico para ofrecer diversos beneficios. Se utiliza principalmente en sistemas que requieren el uso de interfaces de usuario, aunque se puede utilizar en muchos tipos de aplicaciones. Se basa en la separación del código en tres capas distintas según la funcionalidad, Modelos, Vistas y Controladores [14].

El modelo es el backend que tiene la lógica de datos. [15]

La vista es el frontend o la interfaz gráfica de usuario. [15]

El controlador es el que controla como se muestran los datos de la aplicación. [15]

2.3 Herramientas de Desarrollo

En la **TABLA I** se observan las herramientas utilizadas debido a la estructura del proyecto y los requerimientos.

TABLA I. Herramientas fundamentales del proyecto

Herramienta	Descripción	Justificación
Django	Django es un framework de código abierto escrito en Python de alto nivel con el que se fomenta el desarrollo limpio, práctico y rápido [12].	Se utiliza este framework porque es el mejor y el más práctico según los requerimientos del proyecto y porque se utiliza Python.
Python	Python es un lenguaje de programación muy utilizado en el desarrollo de aplicaciones web, desarrollo de software, machine learning y ciencia de datos [16].	Se utiliza este lenguaje de programación ya que es el más práctico de utilizar para web scraping debido a sus librerías.
MongoDB	Es una base de datos NoSql que tiene una gran escalabilidad, flexibilidad y posee un modelo de indexación avanzado [17].	Se utiliza MongoDB porque uno de los requerimientos del proyecto es utilizar una base de datos NoSQL y esta base de datos en tiempo real es la mejor para lo requerido.
Bs4	Es una biblioteca de Python que sirve para analizar documentos en HTML [18].	La biblioteca de beautiful soup nos permite realizar el scrapeo de las páginas web mediante scripts de python.
Instaloader	Es una herramienta que permite descargar todo tipo de contenido de Instagram como fotos, videos, textos y	Instaloader se utiliza porque es una biblioteca que permite extraer información de Instagram la que se

	otros metadatos de Instagram [19].	necesita para el proyecto como tal.
Github	Es la plataforma web visual en la cual se toma a consideración las características de crear proyectos colaborativos en repositorios y sin dificultad a la presentación, pueden ser accedidos desde lugares diferentes por medio de internet [20].	Visualiza el entorno del proyecto de una manera más gráfica y permite administrar el código por etapas a conveniencia del usuario.

3 RESULTADOS

A continuación, se muestran los resultados de los procedimientos realizados en cada uno de los pasos en el proyecto que son presentados en cada uno de los sprints aplicando la metodología RAD:

Sprint 0. Configuración del ambiente de desarrollo

Este sprint inicial hace referencia a la configuración del entorno de desarrollo, asegurándose de que todo lo necesario para empezar con este proceso y codificación esté en óptimas condiciones y con un correcto funcionamiento. Los resultados abarcados en este sprint son:

- VSCode en funcionamiento
- Proyecto en la base de datos NoSQL (MongoDB)

VSCode en funcionamiento

Se prepara el editor de código fuente creando un nuevo proyecto con el lenguaje de programación que se utiliza para el desarrollo de la aplicación como se visualiza en la **Fig 2.**

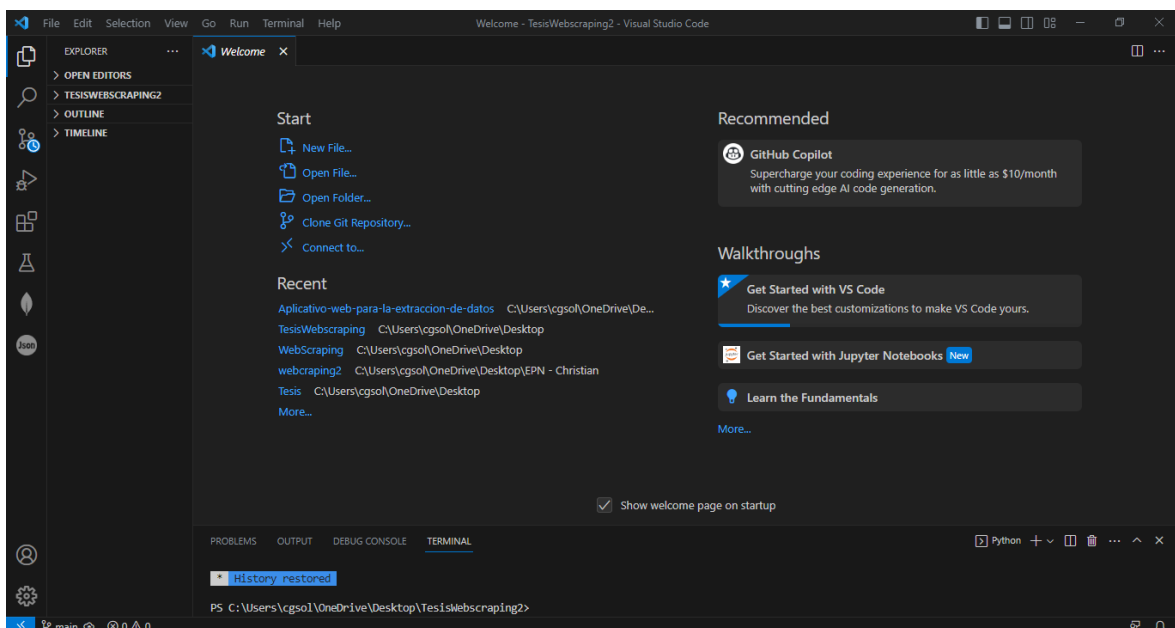


Fig 2. Ejecución del editor de código (VSCode)

Proyecto en la base de datos NoSQL (MongoDB)

Se prepara la base de datos para el proyecto a desarrollar y se crea un proyecto dentro de la misma para conectar con la aplicación como se visualiza en la **Fig 3**.

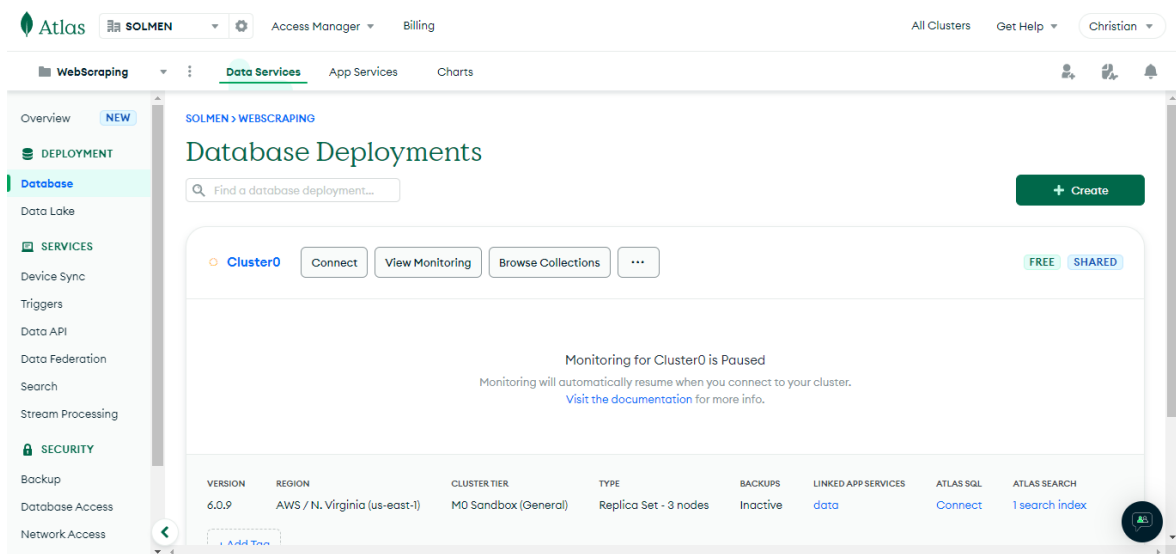


Fig 3. Base de datos en MongoDB

Sprint 1. Elección del lenguaje de programación

El sprint 1 se enfoca en el inicio de la codificación con el lenguaje de programación JavaScript en un inicio y el desarrollo del script para realizar web scraping. Los resultados abarcados en este sprint son:

- Uso de la herramienta playwright.
- Desarrollo del script con JavaScript.
- Error TimeOut.

Uso de la herramienta playwright

Se utiliza la herramienta playwright que permite realizar pruebas end to end en aplicaciones web de manera rápida y práctica [21] para realizar el webscraping mediante Google Chrome como se muestra en la **Fig 4**.

```
JS app.js x {} package-lock.json
JS app.js > getResultFromGoogle
1 import {chromium} from 'playwright';
2
3 // generar resultados de google
4
5
```

Fig 4. Se importa chromium desde playwright

Desarrollo del script con JavaScript

Se utiliza el lenguaje de programación JavaScript para codificar el script con la herramienta playwright como se muestra en la Fig 5.

```
File Edit Selection View Go Run Terminal Help app.js - TesisWebscraping - Visual Studio Code
EXPLORER
> OPEN EDITORS
TESISWEBSCRAPING
  > node_modules
  JS app.js
  {} package-lock.json
  {} package.json
OUTLINE
TIMELINE
JS app.js x {} package-lock.json
JS app.js > getResultFromGoogle
1 import {chromium} from 'playwright';
2
3 // generar resultados de google
4
5
6 async function getResultFromGoogle(query){
7
8   const browser = await chromium.launch();
9   const page = await browser.newPage();
10  await page.goto('https://www.google.com/');
11  await page.waitForSelector('input[name="q"]');
12  await page.type('input[name="q"]', query);
13  await page.keyboard.press('Enter');
14  await page.waitForNavigation({waitUntil: 'networkidle'});
15
16
17
18  const ListadoResultados = await page.evaluate(() => {
19    let resultados = [];
20    document.querySelectorAll('div[data-snhf] div a').forEach((anchor, index) =>{
21      resultados.push({
22        index : index,
23        title : anchor.innerText,
24        url : anchor.href,
25      });
26    });
27    return resultados;
28  });
29  console.log(ListadoResultados);
30  await browser.close();
31
32
```

Fig 5. Script desarrollado con JavaScript

Error TimeOut

Se ejecuta el script de JavaScript con la herramienta playwright y se presenta el error TimeOut 30000ms exceeded que se visualiza en la Fig 6.

```
17
18   const ListadoResultados = await page.evaluate(()=> {
PROBLEMS  OUTPUT  DEBUG CONSOLE  TERMINAL
PS C:\Users\cgsol\OneDrive\Desktop\TesisWebScraping> node app.js
node:internal/process/promises:279
    triggerUncaughtException(err, true /* fromPromise */);
    ^
page.waitForSelector: Timeout 30000ms exceeded.
===== logs =====
waiting for locator('input[name="q"]') to be visible
waiting for locator('input[name="q"]') to be visible
=====
    at getResultFromGoogle (file:///C:/Users/cgsol/OneDrive/Desktop/TesisWebScraping/app.js:11:16) {
  name: 'TimeoutError'
}
PS C:\Users\cgsol\OneDrive\Desktop\TesisWebScraping> node app.js
```

Fig 6. Error TimeOut 30000ms exceeded.

Sprint 2. Cambio del lenguaje de programación

El sprint 2 se enfoca al cambio de lenguaje de programación debido a la practicidad del nuevo lenguaje que es Python, con más variedad de librerías y documentación que permiten que desarrollo del sistema web sea de una manera más rápida, práctica y eficiente. Los resultados abarcados en este sprint son:

- Creación del nuevo proyecto con Python.

Creación del nuevo proyecto con Python.

Se crea el nuevo proyecto con el lenguaje de programación Python para evitar errores suscitados y que puedan suceder a futuro, como se visualiza en la **Fig 7**.

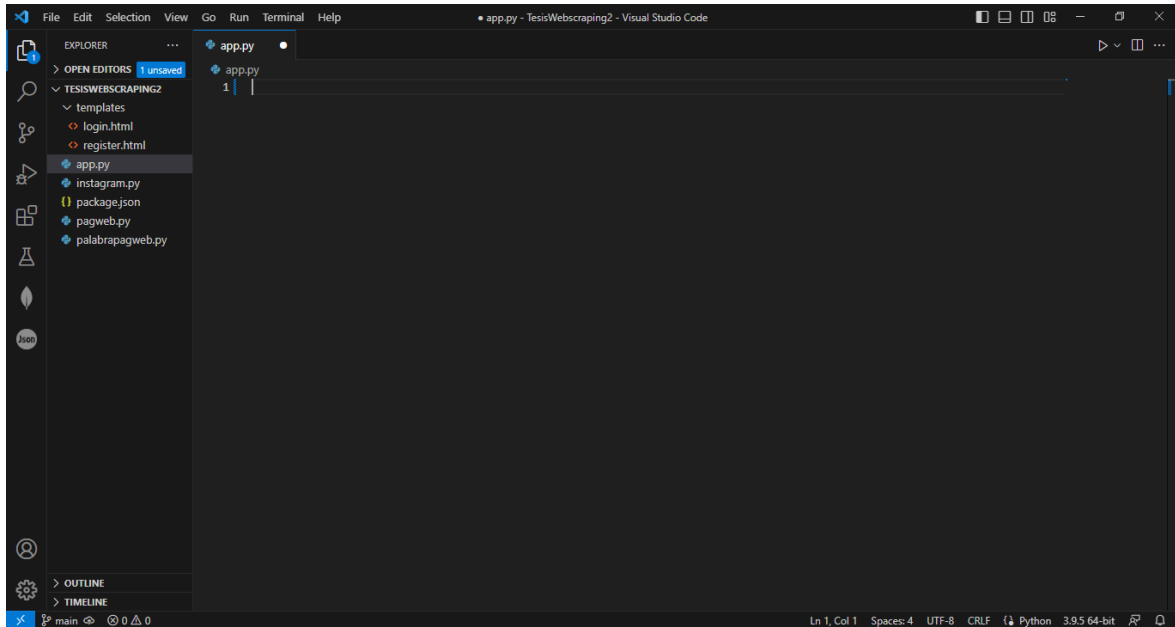


Fig 7. Nuevo proyecto en python

***Sprint 3.* Codificación de Scripts en pyhton para web scraping**

El sprint 3 trata sobre la codificación de los scripts que van a realizar web scraping con el lenguaje de programación Python. Los resultados abarcados en este sprint son:

- Script para extraer toda la información de una página web.
- Script para realizar conteo de palabras en una página web.
- Script para extraer información de perfiles de Instagram.
- Script autenticación de usuario con flask.

Script para extraer toda la información de una página web.

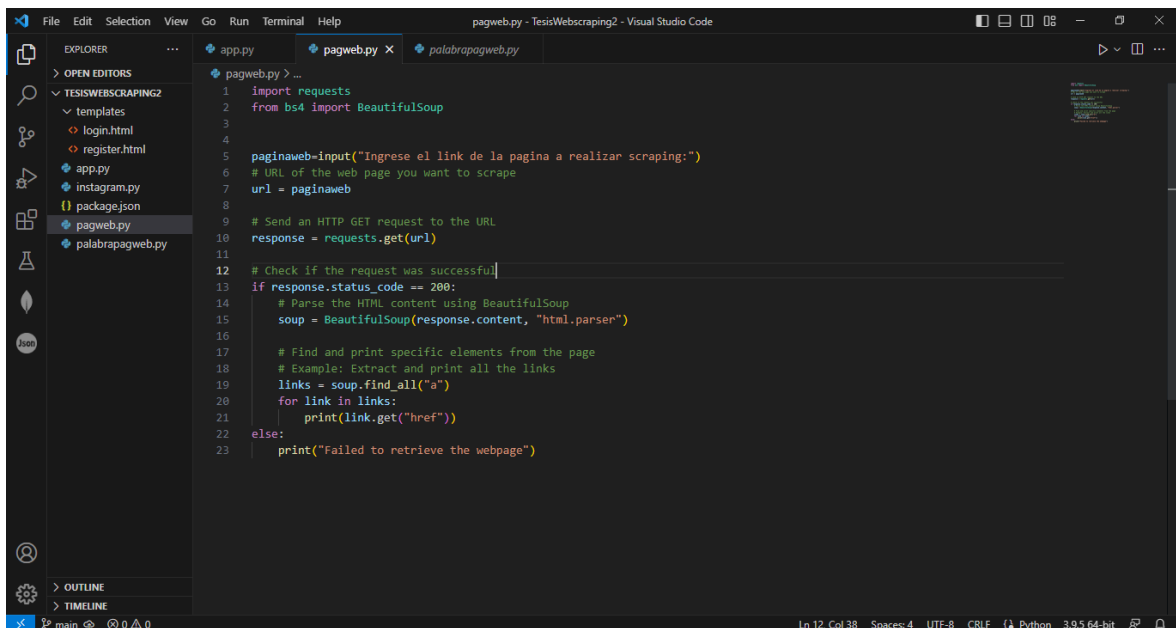
Se desarrolla un script que al ejecutar se debe ingresar el link de una página web y extraer toda la información de esta, utilizando las librerías beautifulsoup4 y requests como se visualiza en las siguientes Fig 8 , Fig 9 y Fig 10.

```
PS C:\Users\cgsol\OneDrive\Desktop\TesisWebScraping2> pip install requests
Requirement already satisfied: requests in c:\users\cgsol\appdata\local\programs\python\python39\lib\site-packages (2.25.1)
Requirement already satisfied: charset<5,>=3.0.2 in c:\users\cgsol\appdata\local\programs\python\python39\lib\site-packages (from requests) (4.0.0)
Requirement already satisfied: idna<3,>=2.5 in c:\users\cgsol\appdata\local\programs\python\python39\lib\site-packages (from requests) (2.10)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\users\cgsol\appdata\local\programs\python\python39\lib\site-packages (from requests) (1.26.6)
Requirement already satisfied: certifi<=2017.4.17 in c:\users\cgsol\appdata\local\programs\python\python39\lib\site-packages (from requests) (2021.5.30)
PS C:\Users\cgsol\OneDrive\Desktop\TesisWebScraping2>
```

Fig 8. Instalación de Requests

```
PS C:\Users\cgsol\OneDrive\Desktop\TesisWebScraping2> pip install beautifulsoup4
Requirement already satisfied: beautifulsoup4 in c:\users\cgsol\appdata\local\programs\python\python39\lib\site-packages (4.9.3)
Requirement already satisfied: soupsieve>1.2 in c:\users\cgsol\appdata\local\programs\python\python39\lib\site-packages (from beautifulsoup4) (2.2.1)
PS C:\Users\cgsol\OneDrive\Desktop\TesisWebScraping2>
```

Fig 9. Instalación de beautifulsoup4



```
pagweb.py > ...
1 import requests
2 from bs4 import BeautifulSoup
3
4
5 paginaweb=input("Ingrese el link de la pagina a realizar scraping:")
6 # URL of the web page you want to scrape
7 url = paginaweb
8
9 # Send an HTTP GET request to the URL
10 response = requests.get(url)
11
12 # Check if the request was successful
13 if response.status_code == 200:
14     # Parse the HTML content using BeautifulSoup
15     soup = BeautifulSoup(response.content, "html.parser")
16
17     # Find and print specific elements from the page
18     # Example: Extract and print all the links
19     links = soup.find_all("a")
20     for link in links:
21         print(link.get("href"))
22 else:
23     print("Failed to retrieve the webpage")
```

Fig 10. Codificación del script para extraer información de una página web

Se ejecuta el script y realiza el scraping ingresando el enlace de la página web de la que se desea extraer la información como se muestra en la Fig 11.

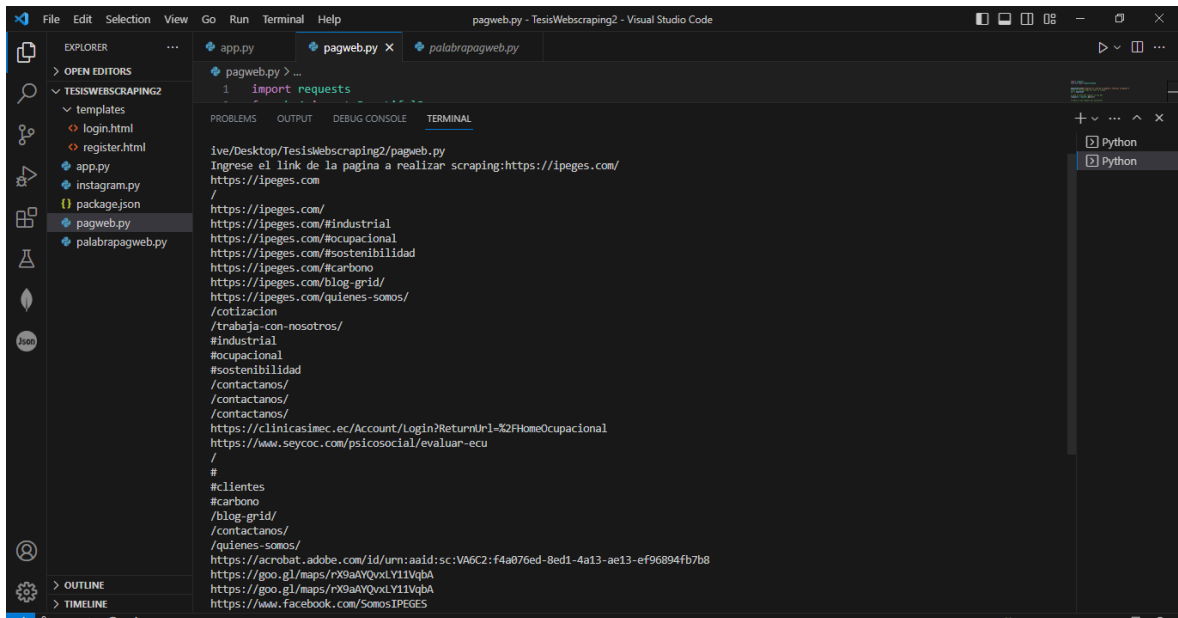


Fig 11. Visualización de la información extraída

Script para realizar conteo de palabras en una página web

Se desarrolla un script que al ejecutar se debe ingresar el enlace de una página web y una palabra que se requiera buscar para realizar un conteo de esta palabra y obtener la información de cuántas veces se la encontró en la página web para que requerimos la instalación de librerías como requests y re lo que se muestra en Fig 12.

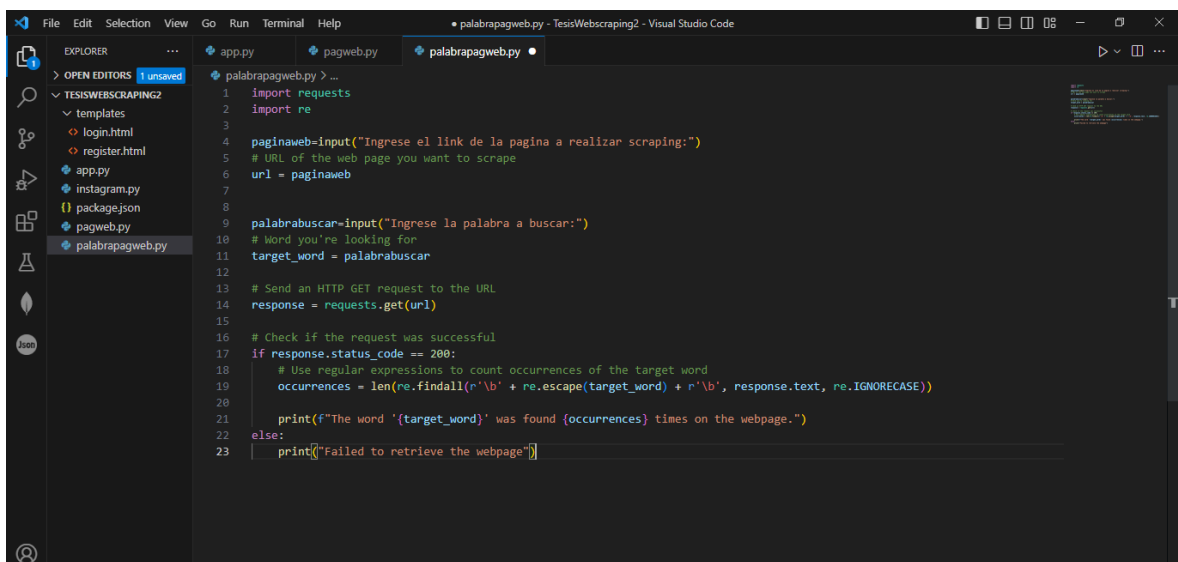


Fig 12. Script para conteo de palabras en una página web.

Se ejecuta el script y al ingresar el enlace de la página web y la palabra de la cual se desea extraer este tipo de información para que realice es scraping y se visualicen los resultados como se muestra en la Fig 13.

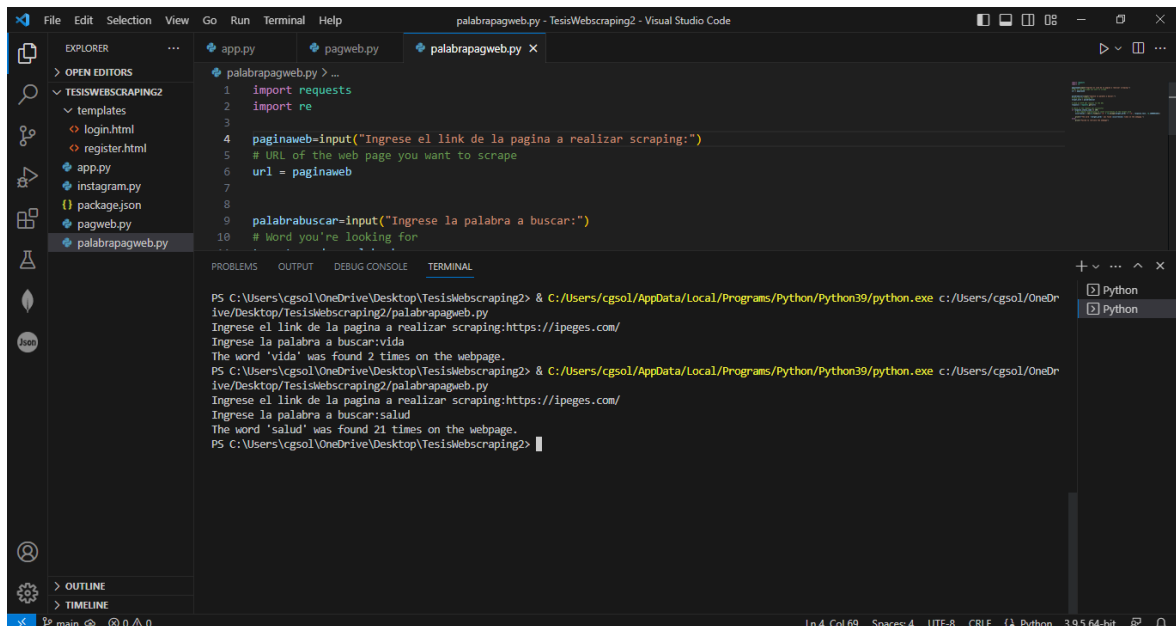


Fig 13. Ejecución y visualización de los datos extraídos del conteo de palabras

Script para extraer información de perfiles de Instagram

Se desarrolla un script utilizando las librerías instalador y pandas para la extracción de datos de perfiles de la red social Instagram y visualice la información como se muestra en Fig 14, Fig 15 y Fig 16.

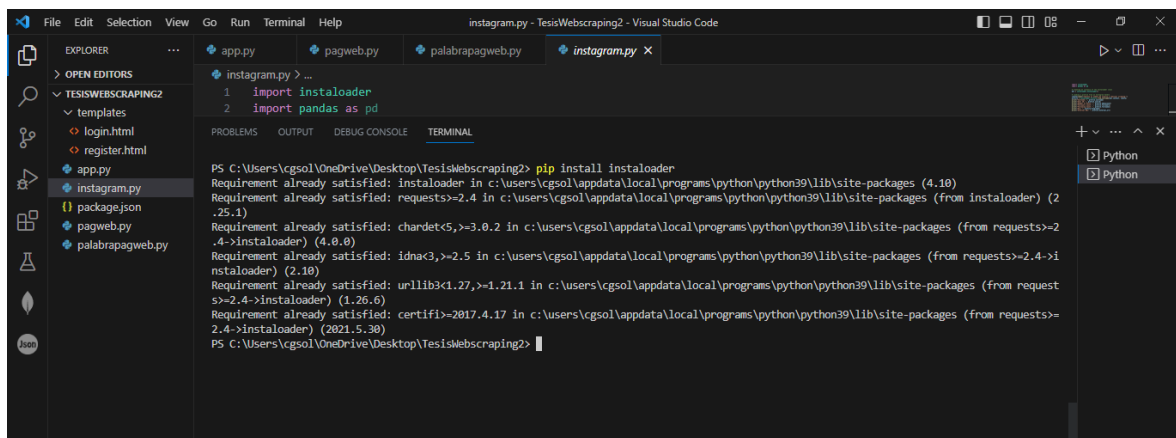


Fig 14. Instalación de Instalador

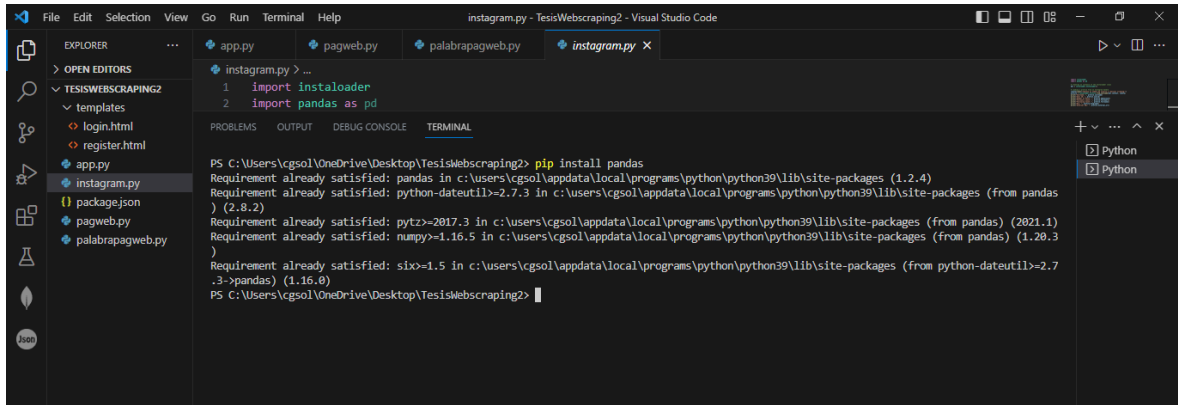


Fig 15. Instalación de pandas

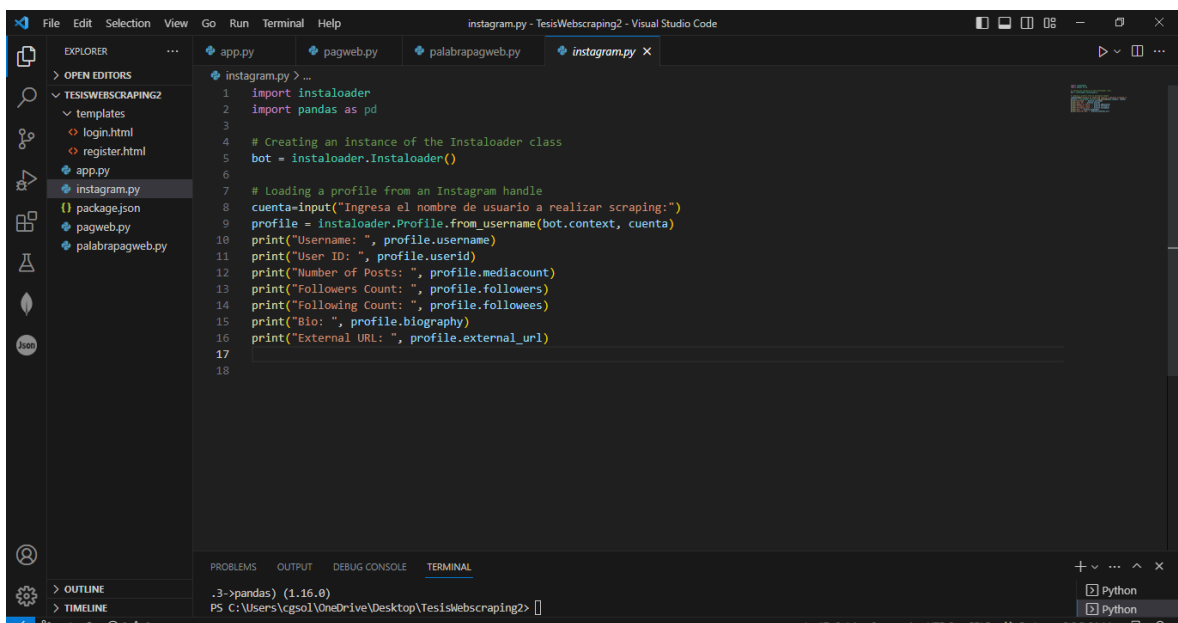


Fig 16. Script para la extracción de datos de Instagram

Se ejecuta el script, se ingresa el usuario de la red social del que se desea extraer la información del perfil para luego visualizarla como se muestra en la Fig 17.

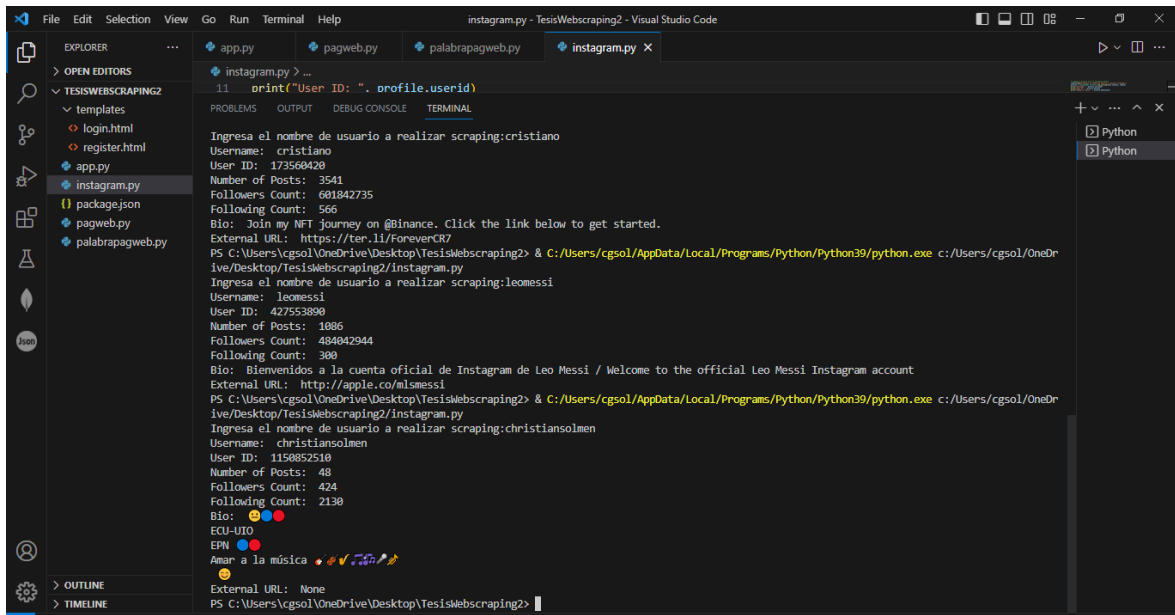


Fig 17. Ejecución y visualización de la información de los perfiles

Script autenticación de usuario con flask.

Se codifica un script para la autenticación de usuario ya conectado a la base de datos NoSQL (MongoDB), realizando la conexión con la base pero generando un error en el script como se muestra en Fig 18, Fig 19 y Fig 20.

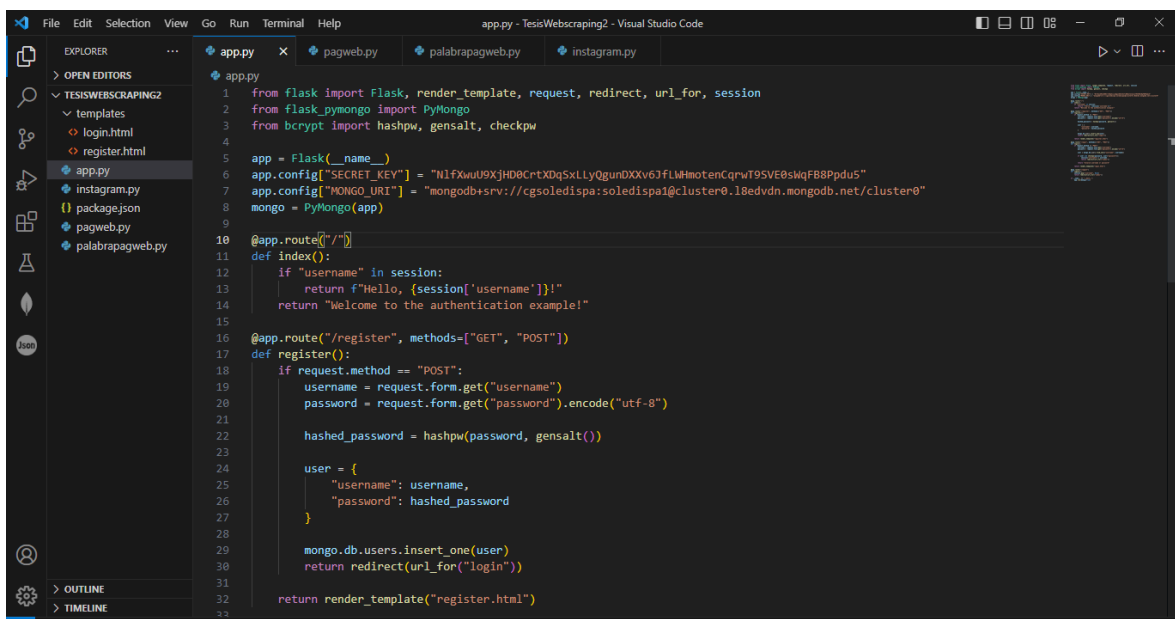


Fig 18. Primera parte del Script de autenticación con flask.

```
28
29     mongo.db.users.insert_one(user)
30     return redirect(url_for("login"))
31
32     return render_template("register.html")
33
34 @app.route("/login", methods=["GET", "POST"])
35 def login():
36     if request.method == "POST":
37         username = request.form.get("username")
38         password = request.form.get("password").encode("utf-8")
39
40         user = mongo.db.users.find_one({"username": username})
41
42         if user and checkpw(password, user["password"]):
43             session["username"] = username
44             return redirect(url_for("index"))
45
46         return "Invalid username or password"
47
48     return render_template("login.html")
49
50 @app.route("/logout")
51 def logout():
52     session.pop("username", None)
53     return redirect(url_for("index"))
54
55 if __name__ == "__main__":
56     app.run(debug=True)
```

Fig 19. Segunda parte del Script de autenticación con flask.

Se ejecuta el script y se realiza la conexión con la base de datos pero el editor de código muestra un error como se visualiza en Fig 20 y Fig 21.

```
External URL: None
PS C:\Users\cgsol\OneDrive\Desktop\TesisWebScraping2> & C:/Users/cgsol/AppData/Local/Programs/Python/Python39/python.exe c:/Users/cgsol/OneDrive/Desktop/TesisWebScraping2/app.py
C:\Users\cgsol\AppData\Local\Programs\Python\Python39\lib\site-packages\cryptography\x509\base.py:583: CryptographyDeprecationWarning: Parsed a negative serial number, which is disallowed by RFC 5280.
  return rust_x509.load_pem_x509_certificate(data)
* Serving Flask app "app"
* Debug mode: on
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on http://127.0.0.1:5000
Press CTRL+C to quit
* Restarting with stat
C:\Users\cgsol\AppData\Local\Programs\Python\Python39\lib\site-packages\cryptography\x509\base.py:583: CryptographyDeprecationWarning: Parsed a negative serial number, which is disallowed by RFC 5280.
  return rust_x509.load_pem_x509_certificate(data)
* Debugger is active!
* Debugger PIN: 500-536-104
```

Fig 20. Error en la ejecución del script de autenticación con flask

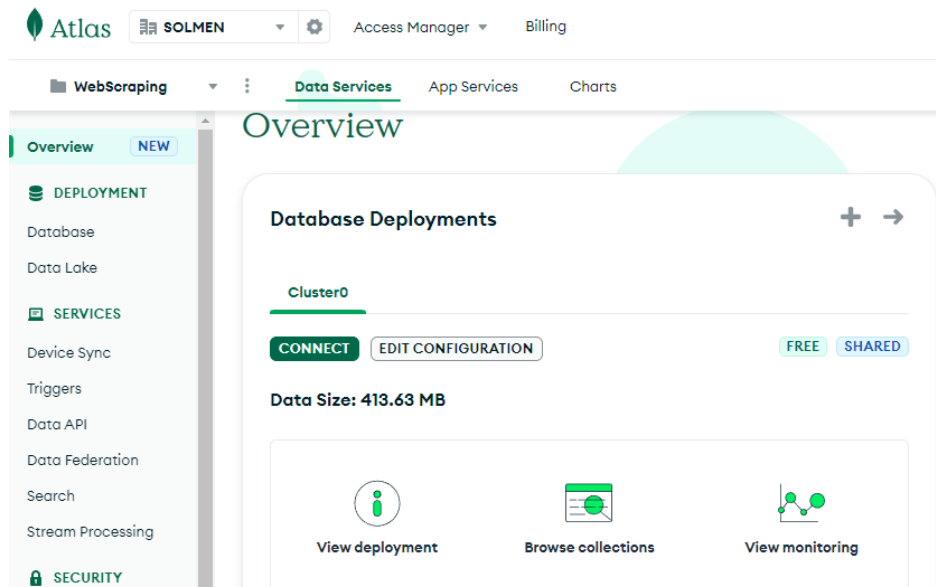


Fig 21. Conexión exitosa con MongoDB

Sprint 4. Implementación de registro de usuarios

El sprint 4 se enfoca en la implementación del registro de usuario para lo cual se utiliza Django ya que con este framework el desarrollo de la autenticación de usuarios se vuelve más práctica al tener facilidad de crear una interfaz amigable que se ajusta a lo requerimientos del sistema de web scraping. Los resultados abarcados en este sprint son:

- Registro de un usuario

Registro de un usuario

El script realiza el proceso del registro de un usuario mediante el host en formato JSON para que los mismos puedan ser almacenados de una manera ordenada y estructurada como se muestra en la Fig 22.

```

request > req.rest > POST //api/auths
4 references
@baseUrl = http://127.0.0.1:8000/
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
###
Send Request
GET {{baseUrl}}api/auths HTTP/1.1
content-type: application/json
###
Send Request
POST {{baseUrl}}/api/auths HTTP/1.1
content-type: application/json
{
  "username": "7chrisred",
  "email": "7chrisred@gmail.com",
  "password": "7chrisred17."
}
###
Send Request
PUT {{baseUrl}}/api/auths/1 HTTP/1.1
content-type: application/json
{
  "username": "Chris Redfield",
  "email": "chris-redfield@gmail.com",
  "password": "chrisred17."
}
###
Send Request
DELETE {{baseUrl}}/api/auths/3 HTTP/1.1
content-type: application/json

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
HTTP/1.1 201 Created
Date: Mon, 21 Aug 2023 07:10:33 GMT
Server: WSGIServer/0.2 CPython/3.10.12
Content-Type: application/json
Vary: Accept, Cookie, origin
Allow: POST, GET, DELETE, OPTIONS
X-Frame-Options: DENY
Content-Length: 117
X-Content-Type-Options: nosniff
Referrer-Policy: same-origin
Cross-Origin-Opener-Policy: same-origin
{
  "id": 8,
  "username": "7chrisred",
  "email": "7chrisred@gmail.com",
  "password": "7chrisred17.",
  "is_superuser": false
}

```

Fig 22. Script del registro de usuarios

Sprint 5. Implementación de Lista de usuarios

El sprint 5 se enfoca en la creación de la lista de usuarios que hayan sido registrados mediante el host y con la estructura definida. Los resultados abarcados en este sprint son:

- Lista de usuarios.

Lista de usuarios

Con el script se genera una lista de todos los usuarios que son registrados, esta lista tiene la estructura de username, email y password como se visualiza en la Fig 23.

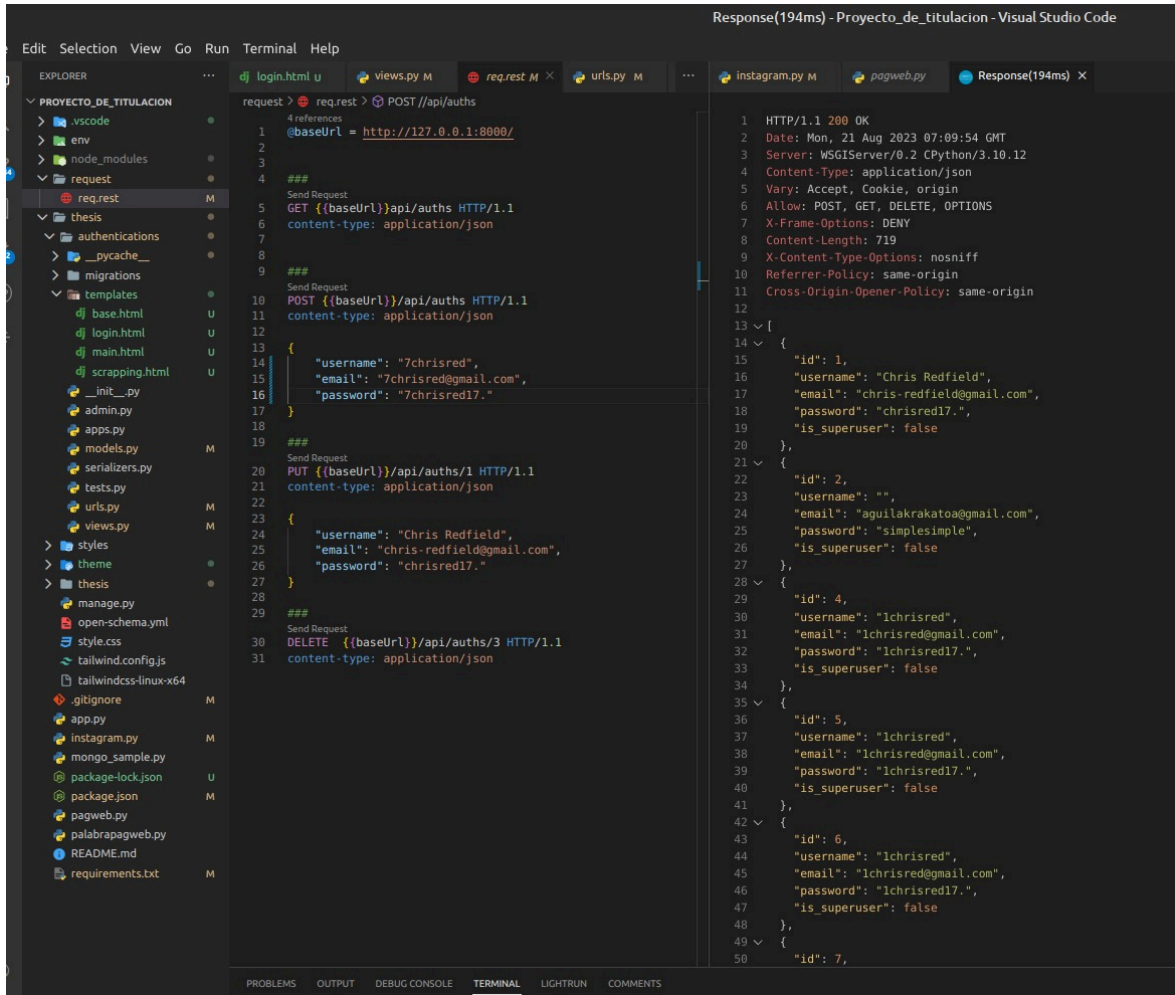


Fig 23. Lista de usuarios

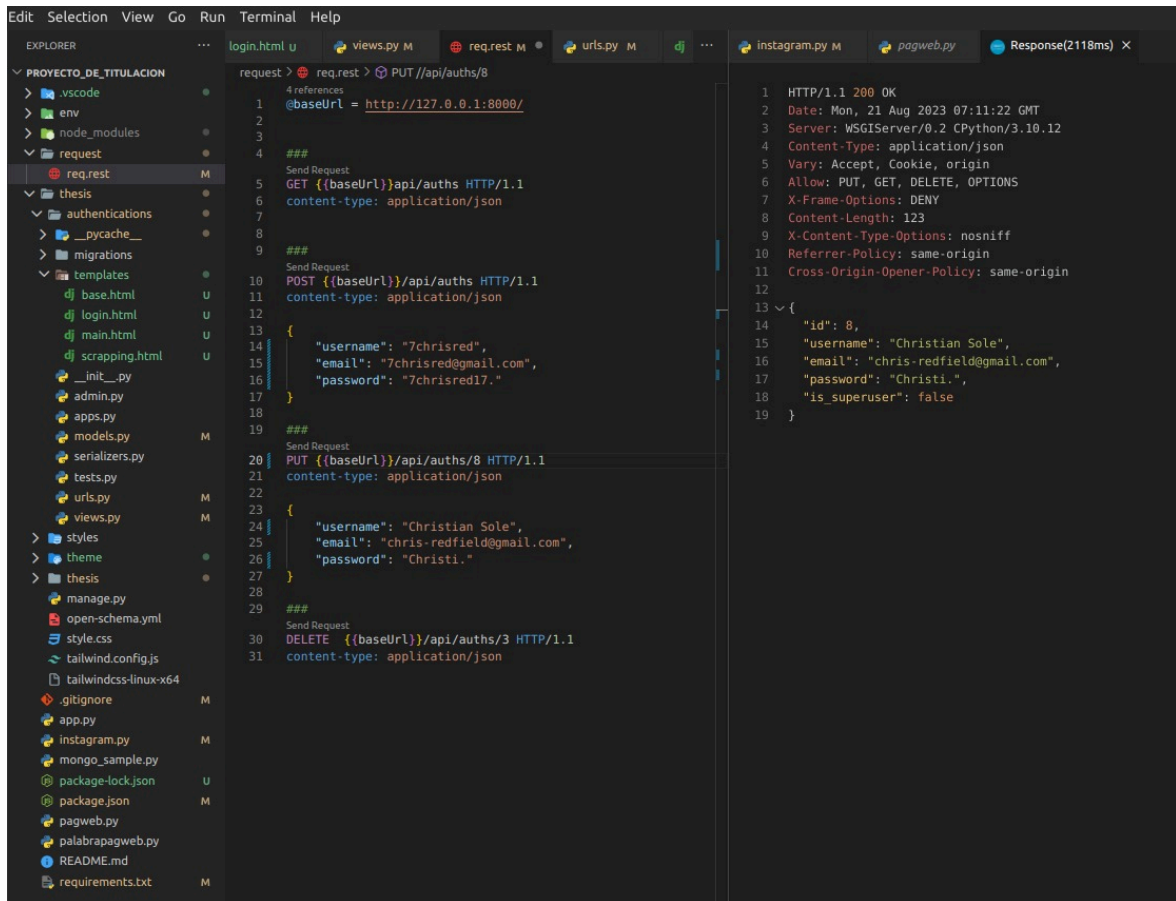
Sprint 6. Implementación de Actualización de usuarios

El sprint 6 se enfoca en la implementación de la función para actualizar un usuario registrado con nueva información en las listas de usuarios registrados sin tener problemas en la autenticación. Los resultados abarcados en este sprint son:

- Actualización de usuarios.

Actualización de usuarios.

Toda la información de un usuario puede ser actualizada mediante el script con Django, esto permite realizar la autenticación luego de la actualización sin tener problemas como se muestra en la Fig 24.



```
request > req.rest > PUT //api/auths/8
4 references
@baseUrl = http://127.0.0.1:8000/
1
2
3
4 ###
5 Send Request
6 GET {{baseUrl}}api/auths HTTP/1.1
7 content-type: application/json
8
9 ###
10 Send Request
11 POST {{baseUrl}}api/auths HTTP/1.1
12 content-type: application/json
13
14 {
15   "username": "7chrisred",
16   "email": "7chrisred@gmail.com",
17   "password": "7chrisred17."
18 }
19 ###
20 Send Request
21 PUT {{baseUrl}}api/auths/8 HTTP/1.1
22 content-type: application/json
23
24 {
25   "username": "Christian Sole",
26   "email": "chris-redfield@gmail.com",
27   "password": "Christi."
28 }
29 ###
30 Send Request
31 DELETE {{baseUrl}}api/auths/3 HTTP/1.1
content-type: application/json

1 HTTP/1.1 200 OK
2 Date: Mon, 21 Aug 2023 07:11:22 GMT
3 Server: WSGIServer/0.2 CPython/3.10.12
4 Content-Type: application/json
5 Vary: Accept, Cookie, origin
6 Allow: PUT, GET, DELETE, OPTIONS
7 X-Frame-Options: DENY
8 Content-Length: 123
9 X-Content-Type-Options: nosniff
10 Referrer-Policy: same-origin
11 Cross-Origin-Opener-Policy: same-origin
12
13 {
14   "id": 8,
15   "username": "Christian Sole",
16   "email": "chris-redfield@gmail.com",
17   "password": "Christi.",
18   "is_superuser": false
19 }
```

Fig 24. Actualización de usuarios

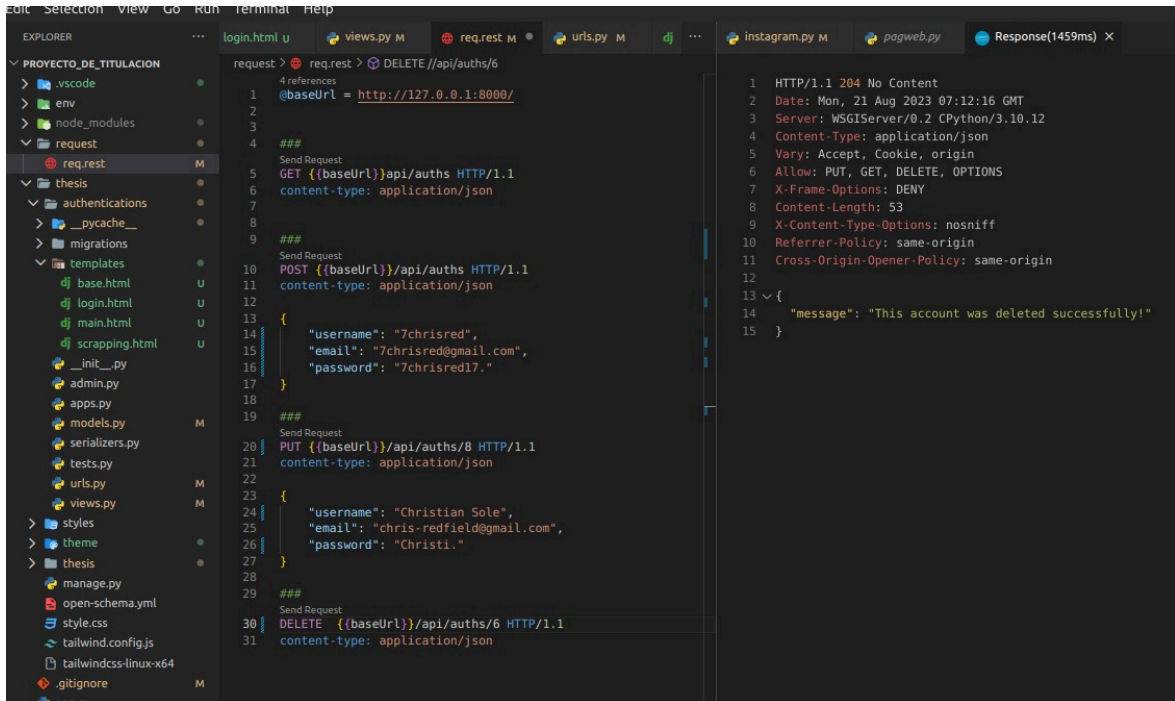
Sprint 7. Implementación de Borrar usuarios

En el sprint 7 se desarrolla la función que permite eliminar un usuario registrado y que no conste en la base, así cuando se realiza el proceso de autenticación aparecerá como usuario inexistente. Los resultados que abarca este sprint son:

- Borrar usuarios.

Borrar usuarios

En el script se crea una función que elimina el usuario que previamente debe estar registrado y al autenticar muestra el mensaje de que el usuario no existe como se muestra en la Fig 25.



```
request > req.rest > DELETE //api/auths/6
1 @baseUrl = http://127.0.0.1:8000/
2
3
4 ###
5 Send Request
6 GET {{baseUrl}}api/auths HTTP/1.1
7 content-type: application/json
8
9 ###
10 Send Request
11 POST {{baseUrl}}/api/auths HTTP/1.1
12 content-type: application/json
13
14 {
15   "username": "7chrisred",
16   "email": "7chrisred@gmail.com",
17   "password": "7chrisred17."
18 }
19 ###
20 Send Request
21 PUT {{baseUrl}}/api/auths/8 HTTP/1.1
22 content-type: application/json
23
24 {
25   "username": "Christian Sole",
26   "email": "chris-redfield@gmail.com",
27   "password": "Christi."
28 }
29 ###
30 Send Request
31 DELETE {{baseUrl}}/api/auths/6 HTTP/1.1
content-type: application/json

1 HTTP/1.1 204 No Content
2 Date: Mon, 21 Aug 2023 07:12:16 GMT
3 Server: WSGIServer/0.2 CPython/3.10.12
4 Content-Type: application/json
5 Vary: Accept, Cookie, origin
6 Allow: PUT, GET, DELETE, OPTIONS
7 X-Frame-Options: DENY
8 Content-Length: 53
9 X-Content-Type-Options: nosniff
10 Referrer-Policy: same-origin
11 Cross-Origin-Opener-Policy: same-origin
12
13 {
14   "message": "This account was deleted successfully!"
15 }
```

Fig 25. Borrar usuarios

Sprint 8. Implementación Login para el usuario en la web

El sprint 8 se enfoca en el desarrollo de la interfaz en la web con el uso de Django para el login del usuario. Los resultados que abarca este sprint son:

- Login para el usuario en la web

Login para el usuario en la web

En el script se desarrolla el código en html donde el usuario hace login mediante la interfaz proporcionada igual por Django como se muestra en la Fig 26.

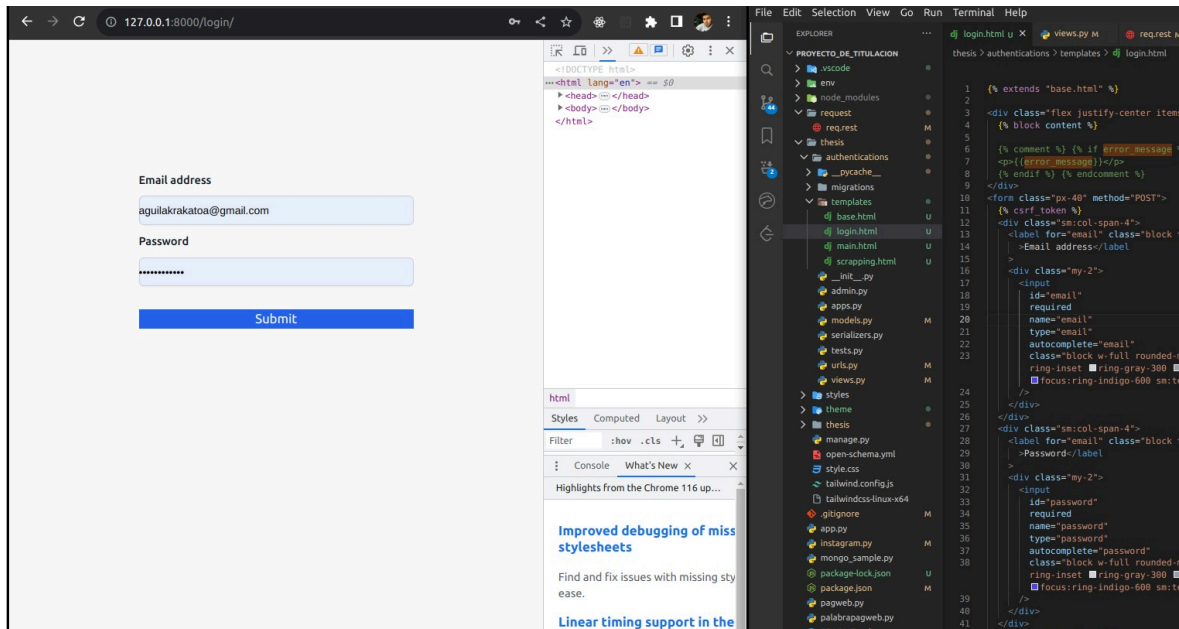


Fig 26. Login para el usuario

Sprint 9. Implementación de scraping en la interfaz web

El sprint 9 se enfoca en el desarrollo de la interfaz web mediante Django y html donde el usuario ingresa y busca tanto usuario de Instagram como página web a scrapear. Los resultados que abarca este sprint son:

- Scraping en la web.

Scraping en la web

En el script se desarrolla un código en html con Django que presenta al usuario una interfaz donde puede ingresar la información para lo que se quiere realizar scraping como se muestra en la Fig 27.

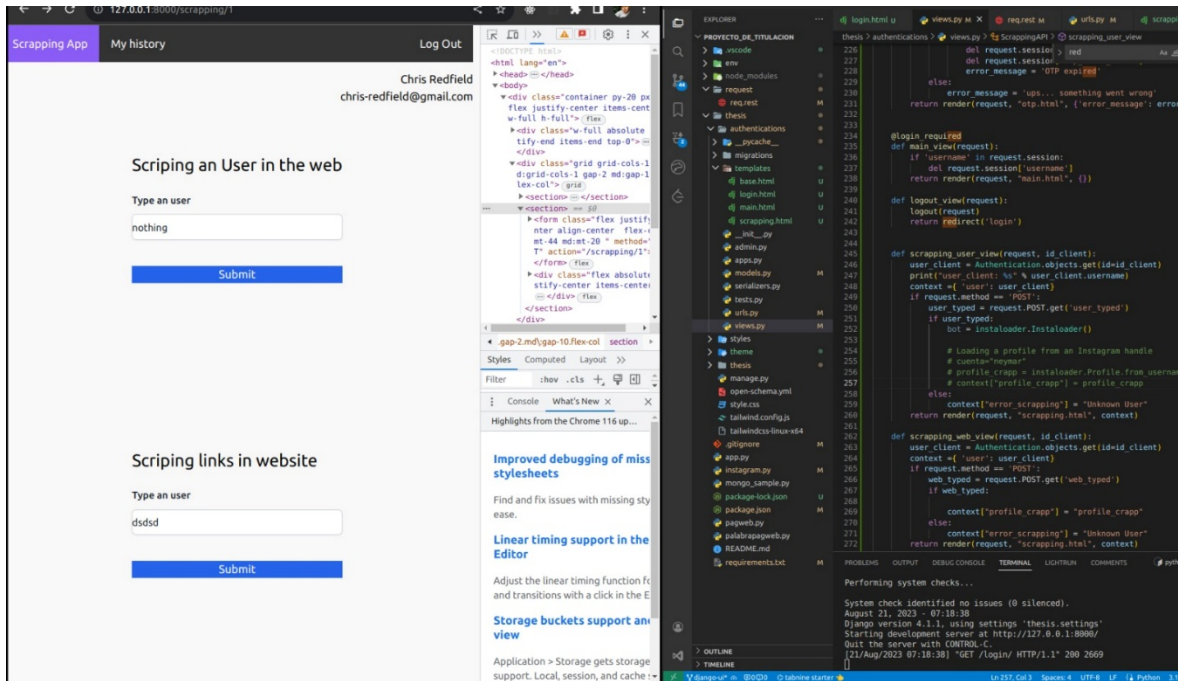


Fig 27. Scraping en la web

Sprint 10. Almacenamiento de usuarios registrados en la base de datos

El sprint 10 se enfoca en almacenar los datos de los usuarios registrados, en la base de datos MongoDB. Los resultados que abarca este sprint son:

- Almacenamiento de usuarios en la base de datos.

Almacenamiento de usuarios en la base de datos

Todos los datos de los usuarios registrados son almacenados en la base de datos NoSQL como se muestra en la Fig 28.

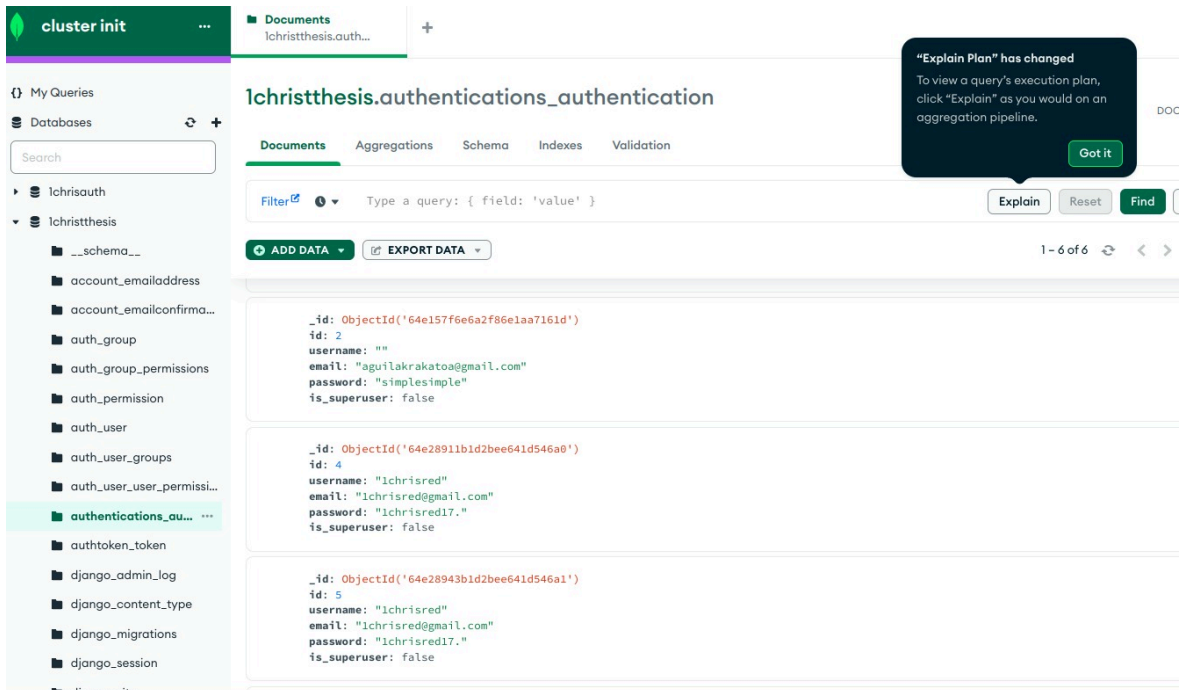


Fig 28. Información en la base de datos

Sprint 11. Scraping anchors en la web

El sprint 11 se enfoca en mostrar los resultados del scraping realizado en la interfaz, con una consulta de la API el servidor arroja los resultados de un sitio web. Los resultados sobre este sprint son:

- Scraping anchors en el sitio web.
- Consulta de la API.

Scraping anchors en el sitio web

Los resultados de realizar scraping se visualizan en la interfaz luego de realizar el login como se muestra en la Fig 28.

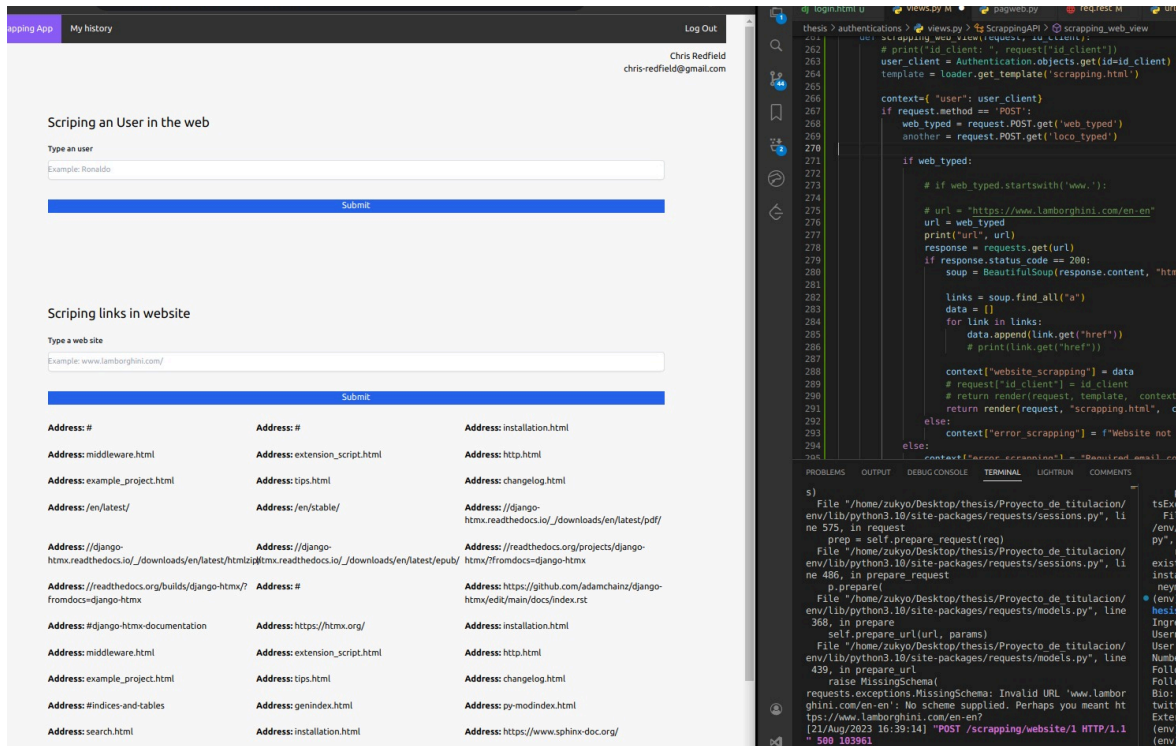


Fig 28. Scraping anchors en la web

Consulta de la API

La consulta de la API al servidor arroja los resultados de anchors de un sitio web como se muestra en la Fig 29.

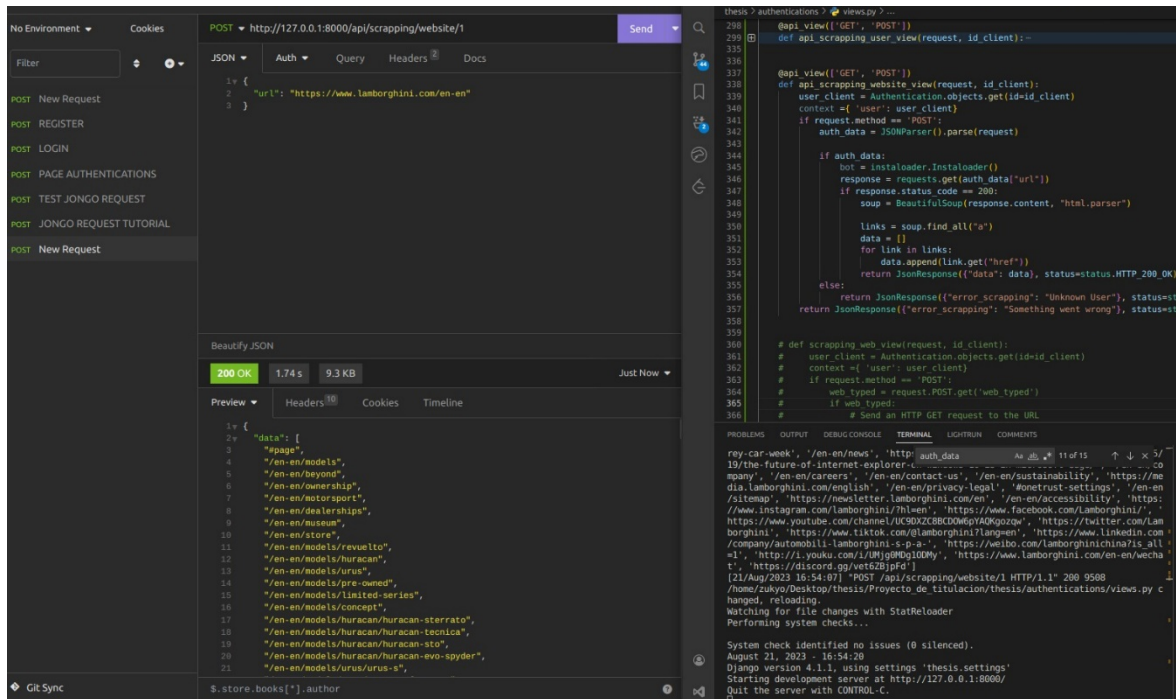


Fig 29. Consulta de la API

Sprint 12. Rutas de Django y Modelos de usuario

El sprint 12 se enfoca en la creación de las rutas con Django y la creación de los modelos de usuarios. Los resultados que abarcan este sprint son:

- Rutas de Django
- Modelos de usuario

Rutas Django

El script muestra las rutas creadas para los datos tanto de autenticación como los datos extraídos de las páginas web como se muestra en las Fig 30 y Fig 31.


```

11
12 from rest_framework_simplejwt.views import (
13     TokenObtainPairView,
14     TokenRefreshView,
15 )
16
17 from drf_spectacular.views import SpectacularAPIView, SpectacularSwaggerView
18
19
20 from .views import login_view, otp_view, main_view, logout_view, scrapping_view
21 # from allauth.account.views import LoginView
22 # from allauth.account.views import SignupView
23 # from allauth.account.views import LogoutView
24
25 urlpatterns = [
26     path("admin/", admin.site.urls),
27     path('', include('authentications.urls')),
28     # path('api/schema', SpectacularAPIView.as_view(), name="schema"),
29     # path('api/schema/docs/', SpectacularSwaggerView.as_view(url_name="schema"))
30
31
32     # path('', main_view, name="main"),
33     path("__reload__/", include("django_browser_reload.urls")),
34
35

```

Fig 30. Definición de las rutas

```

7 # ]
8
9
10 from django.urls import path
11 from .views import AuthenticationAPI, ScrappingAPI
12
13 urlpatterns = [
14     path('api/auths', AuthenticationAPI.authentication_list),
15     path('api/auths/<int:pk>', AuthenticationAPI.authentication_detail),
16     path('api/auths/superuser', AuthenticationAPI.authentication_list_is_superuser),
17     path('login/', ScrappingAPI.login_view, name="login"),
18     path('signup/', ScrappingAPI.signup_view, name="signup"),
19     path('logout/', ScrappingAPI.logout_view, name="logout"),
20     path('scrapping/<id_client>', ScrappingAPI.scrapping_user_view, name="scrapping"),
21     path('scrapping/website/<id_client>', ScrappingAPI.scrapping_web_view, name="website_scrapping"),
22     path('api/scrapping/<id_client>', ScrappingAPI.api_scrapping_user_view, name="api_scrapping"),
23     path('api/scrapping/website/<id_client>', ScrappingAPI.api_scrapping_website_view, name="api_website_scrapping"),
24     path('', ScrappingAPI.main_view, name='main'),
25 ]

```

Fig 31. Rutas Django

Modelos de usuario

Se estructuran los modelos de usuarios mediante djongo y tener la estructura deseada como se muestra en la Fig 32.

```

1 from django.db import models
2 from django import forms
3
4
5
6 class Authentication(models.Model):
7     username = models.CharField(max_length=70, blank=True, default='')
8     email = models.CharField(max_length=200, blank=False)
9     password = models.CharField(max_length=200, blank=False)
10    is_superuser = models.BooleanField(default=False)
11    history = models.ArrayReferenceField(
12        to=models.ObjectIdField(),
13        blank=True,
14        null=True
15    )
16

```

Fig 32. Modelos de usuarios

Sprint 13. Implementación de vistas

El sprint 13 se enfoca en la implementación de vistas de frontend tanto para el login y singup como para el scraping. Los resultados de este sprint son los siguientes:

- Vistas frontend de login y singup
- Vistas frontend del scraping

Vistas Frontend de login y singup

El script contiene las vistas en la interfaz login y singup como se visualiza en la Fig 33.

```

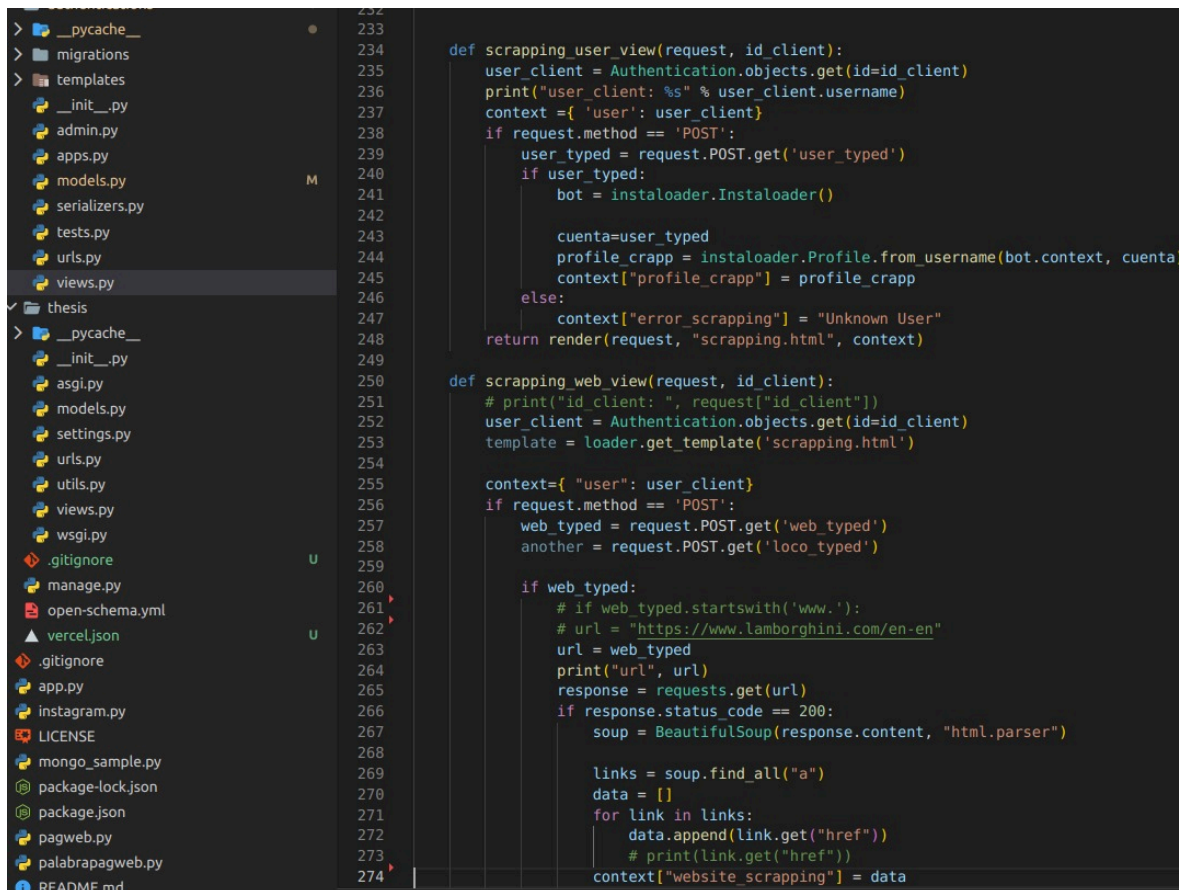
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000

```

Fig 33. Vistas Login y Singup

Vistas Frontend del scraping

El script contiene las vistas en de la interfaz para el web scraping como se visualiza en la Fig 34.



```
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274

def scrapping_user_view(request, id_client):
    user_client = Authentication.objects.get(id=id_client)
    print("user_client: %s" % user_client.username)
    context = { 'user': user_client }
    if request.method == 'POST':
        user_typed = request.POST.get('user_typed')
        if user_typed:
            bot = instaloader.Instaloader()

            cuenta=user_typed
            profile_crapp = instaloader.Profile.from_username(bot.context, cuenta)
            context["profile_crapp"] = profile_crapp
        else:
            context["error_scrapping"] = "Unknown User"
    return render(request, "scrapping.html", context)

def scrapping_web_view(request, id_client):
    # print("id_client: ", request["id_client"])
    user_client = Authentication.objects.get(id=id_client)
    template = loader.get_template('scrapping.html')

    context={ "user": user_client }
    if request.method == 'POST':
        web_typed = request.POST.get('web_typed')
        another = request.POST.get('loco_typed')

        if web_typed:
            # if web_typed.startswith('www.'):
            # url = "https://www.lamborghini.com/en-en"
            url = web_typed
            print("url", url)
            response = requests.get(url)
            if response.status_code == 200:
                soup = BeautifulSoup(response.content, "html.parser")

                links = soup.find_all("a")
                data = []
                for link in links:
                    data.append(link.get("href"))
                    # print(link.get("href"))
            context["website_scrapping"] = data
```

Fig 34. Vistas scraping

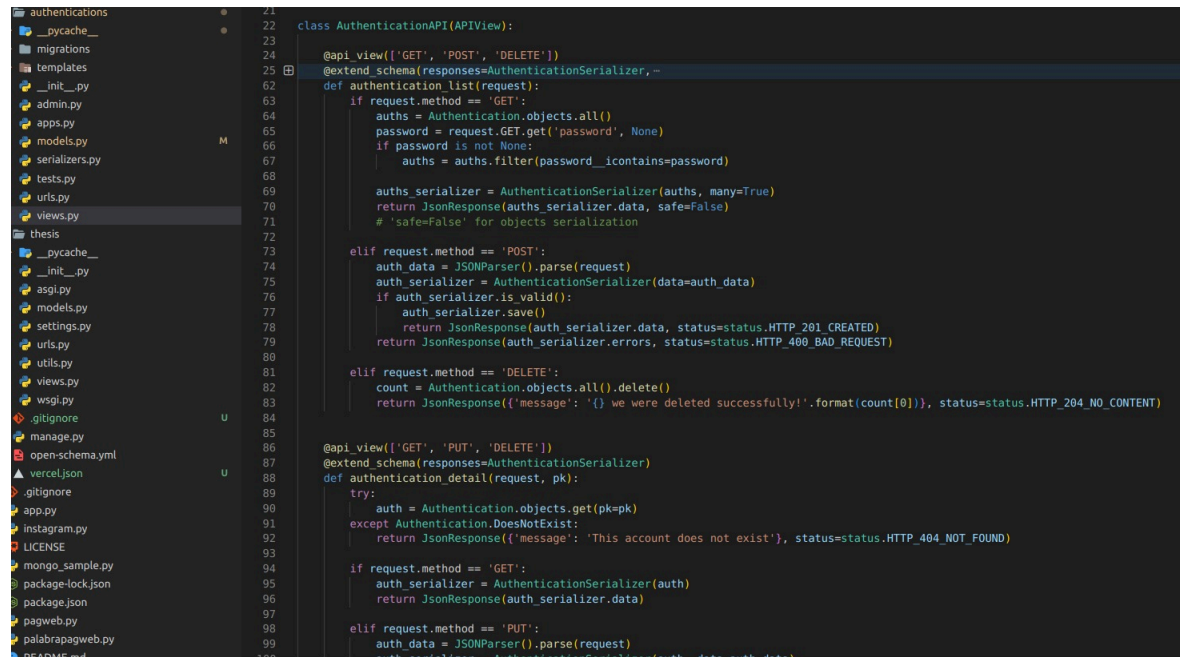
***Sprint 14.* Implementación de API Rest de autenticación y API Rest de listas de usuarios**

En el sprint 14 se desarrolla el consumo de las API Rest tanto de autenticación como de listas de usuarios. Los resultados de este sprint son:

- API Rest de autenticación
- API Rest de Listas de usuarios

API Rest de autenticación

Se codifica la API Rest de autenticación para que arroje los resultados de los usuarios registrados y se realice la autenticación correctamente como se visualiza en la Fig 35.



```
class AuthenticationAPI(APIView):
    @api_view(['GET', 'POST', 'DELETE'])
    @extend_schema(responses=AuthenticationSerializer, ...)
    def authentication_list(request):
        if request.method == 'GET':
            auths = Authentication.objects.all()
            password = request.GET.get('password', None)
            if password is not None:
                auths = auths.filter(password__icontains=password)

            auths_serializer = AuthenticationSerializer(auths, many=True)
            return JsonResponse(auths_serializer.data, safe=False)
            # 'safe=False' for objects serialization

        elif request.method == 'POST':
            auth_data = JSONParser().parse(request)
            auth_serializer = AuthenticationSerializer(data=auth_data)
            if auth_serializer.is_valid():
                auth_serializer.save()
                return JsonResponse(auth_serializer.data, status=status.HTTP_201_CREATED)
            return JsonResponse(auth_serializer.errors, status=status.HTTP_400_BAD_REQUEST)

        elif request.method == 'DELETE':
            count = Authentication.objects.all().delete()
            return JsonResponse({'message': '{} we were deleted successfully!'.format(count[0])}, status=status.HTTP_204_NO_CONTENT)

    @api_view(['GET', 'PUT', 'DELETE'])
    @extend_schema(responses=AuthenticationSerializer)
    def authentication_detail(request, pk):
        try:
            auth = Authentication.objects.get(pk=pk)
        except Authentication.DoesNotExist:
            return JsonResponse({'message': 'This account does not exist'}, status=status.HTTP_404_NOT_FOUND)

        if request.method == 'GET':
            auth_serializer = AuthenticationSerializer(auth)
            return JsonResponse(auth_serializer.data)

        elif request.method == 'PUT':
            auth_data = JSONParser().parse(request)
            auth_serializer = AuthenticationSerializer(auth_data)
            if auth_serializer.is_valid():
                auth_serializer.save()
                return JsonResponse(auth_serializer.data, status=status.HTTP_201_CREATED)
            return JsonResponse(auth_serializer.errors, status=status.HTTP_400_BAD_REQUEST)
```

Fig 35. API Rest de authentication

API Rest de Listas de usuarios

Se codifica la API Rest de Listas de usuarios para que esta arroje los resultados de las listas correctamente y de igual manera no tener problemas con la autenticación de los usuarios como se muestra en la Fig 36.

```
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

Fig 36. API Rest de Listas de Usuarios

Sprint 15. Implementación total del web scraping

En el sprint 15 se implementa todo el sistema para realizar web scraping como la extracción de los perfiles de Instagram, extracción de la información de páginas web, la vista en el fronted y un historial. Los resultados que abarca este sprint son:

- Scraping de perfiles de Instagram.
- Scraping de páginas web.
- Historial de scraping.

Scraping de perfiles de Instagram

Se desarrolla la interfaz mediante Django donde el usuario ingresa el nombre del perfil y se extrae la información de este y la visualiza dentro de la misma interfaz como se muestra en la Fig 37.

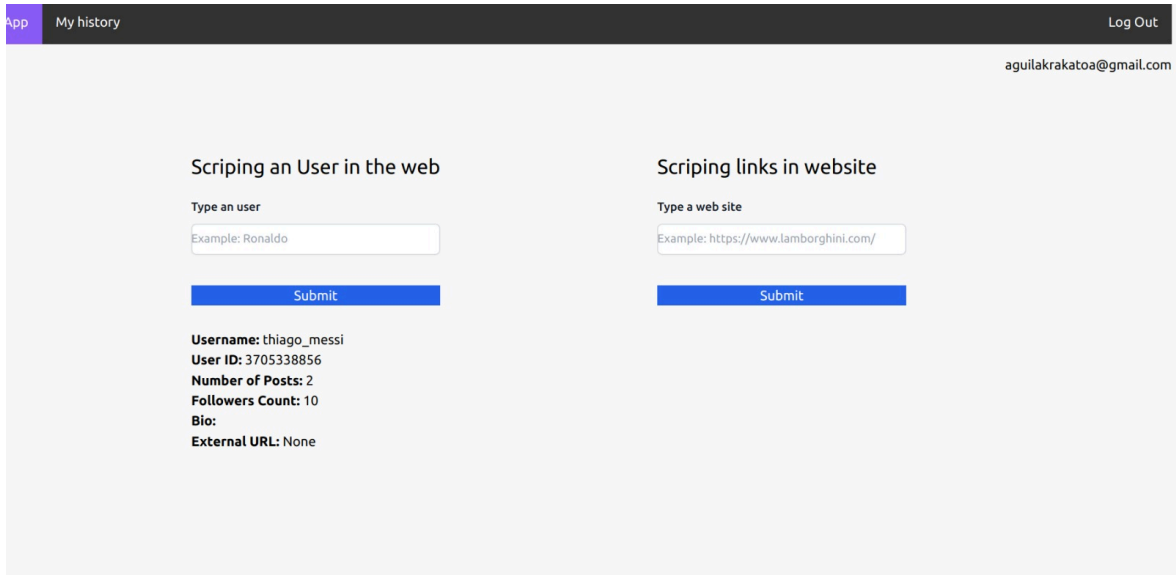


Fig 37. Scraping perfil de usuario

Scraping de páginas web

Se desarrolla la interfaz mediante Django donde el usuario ingresa el enlace de la página web, se extrae toda la información y se visualiza dentro de la interfaz como se muestra en la Fig 38.

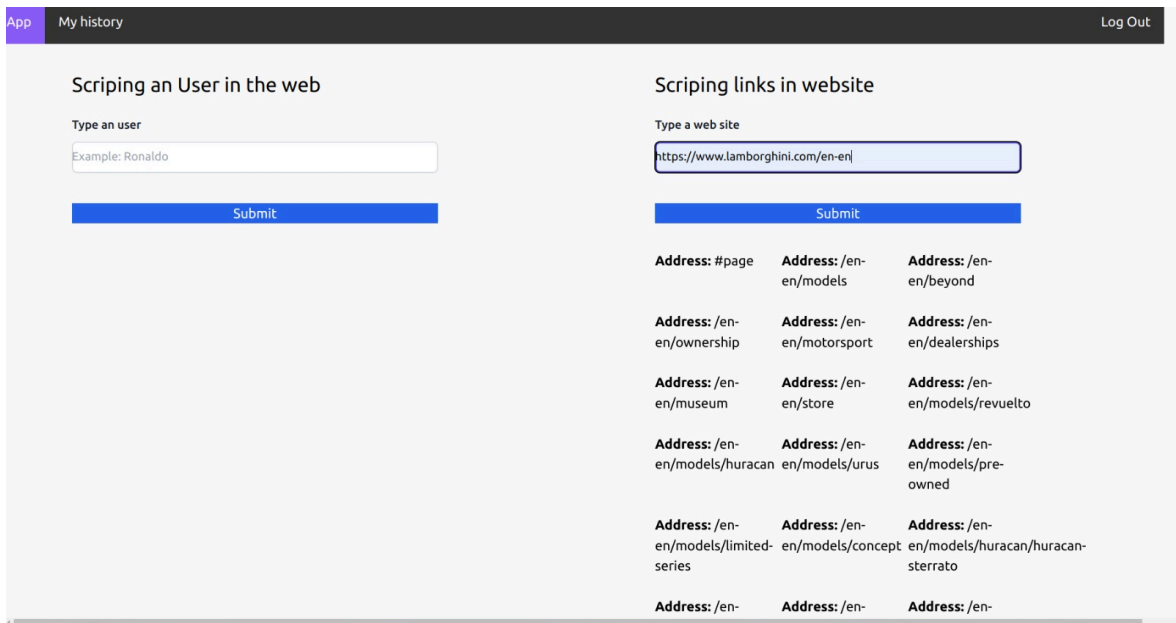
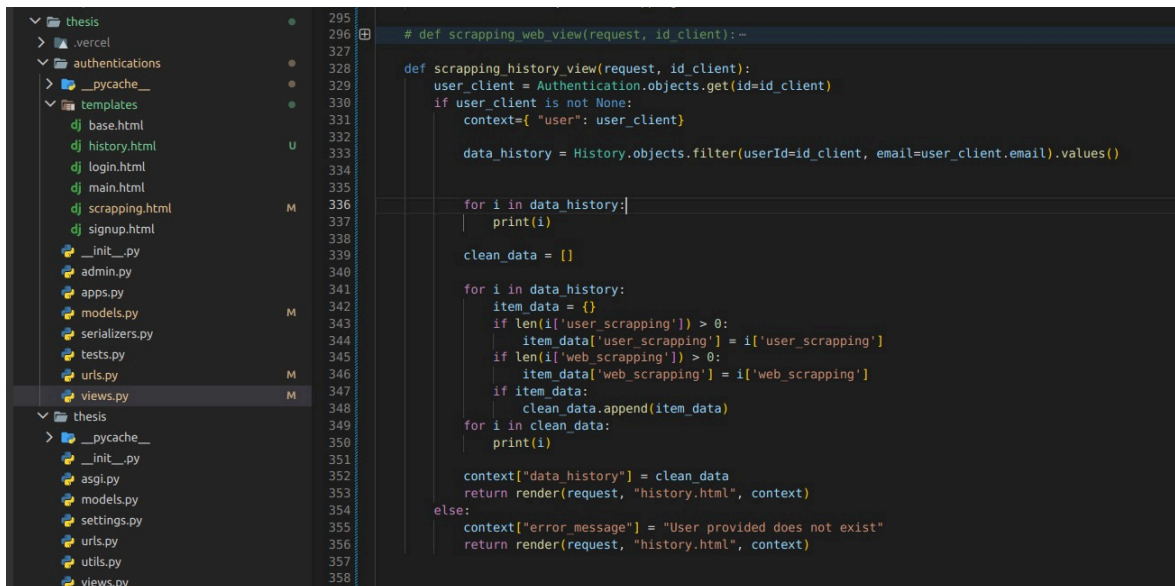


Fig 39. Scriping página web

Historial de web scraping

Se desarrolla la vista del scraping mediante el código y se lo visualiza mediante una pequeña interfaz como se muestra en las Fig 39 y Fig 40.



```
295 # def scrapping_web_view(request, id_client):-
296
327
328 def scrapping_history_view(request, id_client):
329     user_client = Authentication.objects.get(id=id_client)
330     if user_client is not None:
331         context={"user": user_client}
332
333         data_history = History.objects.filter(userId=id_client, email=user_client.email).values()
334
335         for i in data_history:
336             print(i)
337
338         clean_data = []
339
340         for i in data_history:
341             item_data = {}
342             if len(i['user_scraping']) > 0:
343                 item_data['user_scraping'] = i['user_scraping']
344             if len(i['web_scraping']) > 0:
345                 item_data['web_scraping'] = i['web_scraping']
346             if item_data:
347                 clean_data.append(item_data)
348         for i in clean_data:
349             print(i)
350
351         context["data_history"] = clean_data
352         return render(request, "history.html", context)
353     else:
354         context["error_message"] = "User provided does not exist"
355         return render(request, "history.html", context)
356
357
358
```

Fig 39. Vista de frontend del historial

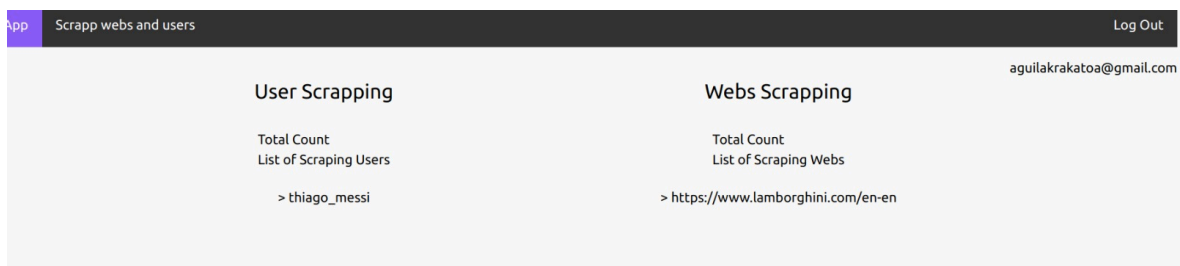


Fig 40. Historial del scraping

Sprint 16. Almacenamiento de los datos del historial en la base de datos

El sprint 16 se enfoca en almacenar los datos extraídos del historial al hacer web scraping en la base de datos MongoDB. Los resultados que abarca este sprint son:

- Base de datos del Historial.

Base de datos del historial

Los datos obtenidos al realizar tanto el web scraping de Instagram como de las páginas web se almacenan en la base de datos de MongoDB como se visualiza en la Fig 41.

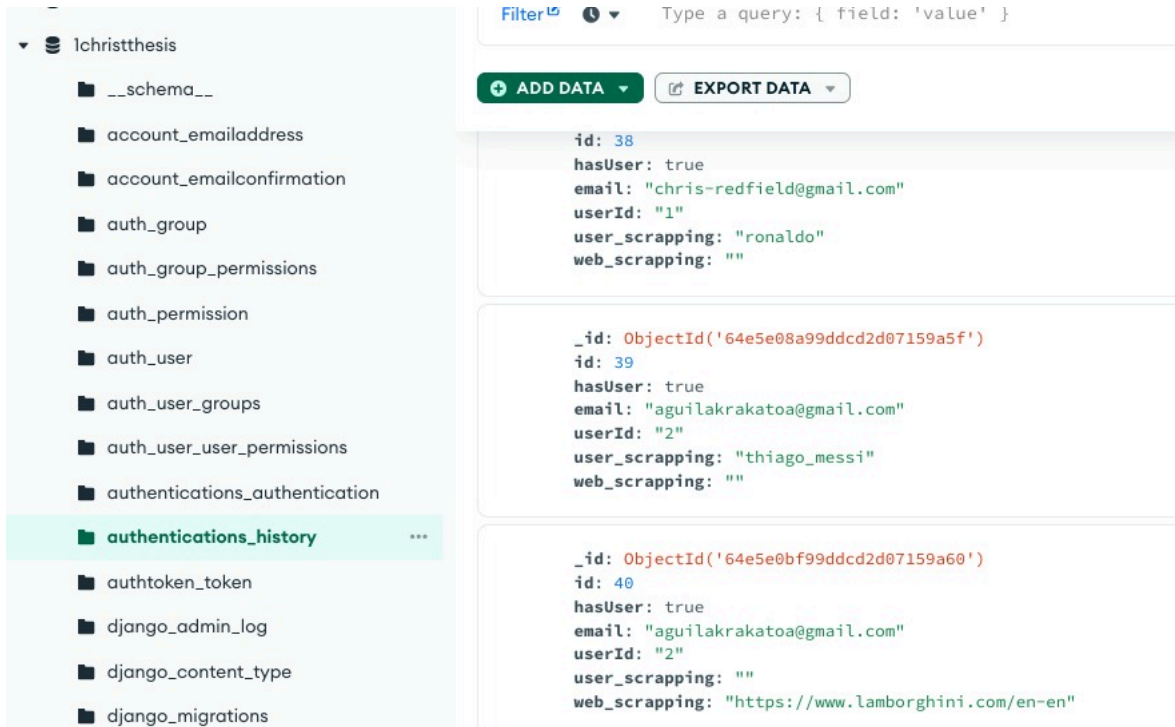


Fig 41. Base de datos del historial de web scraping

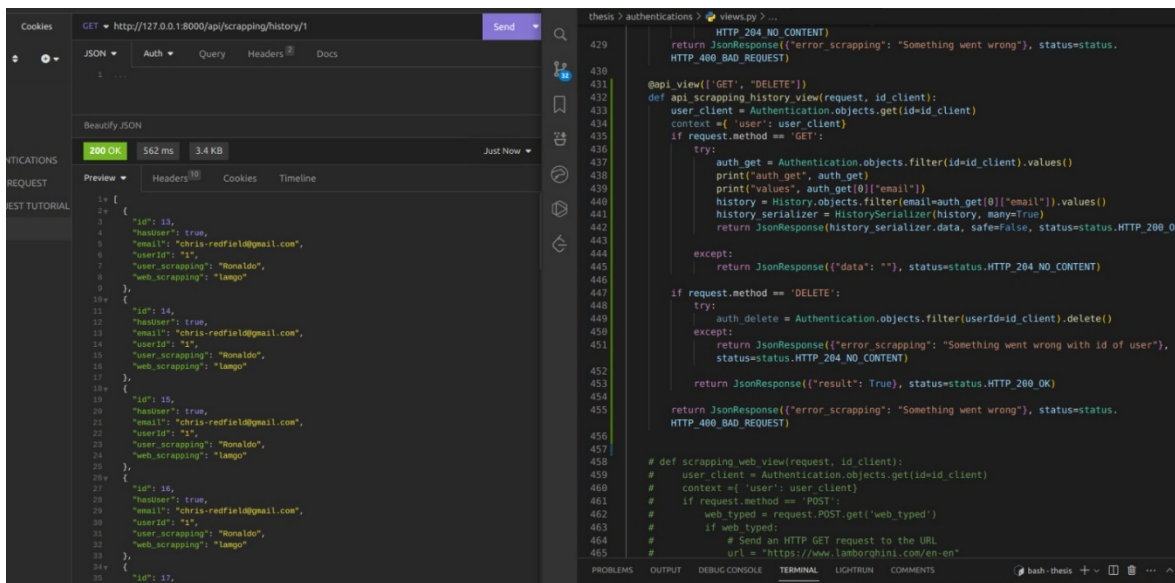


Fig 42. API Rest del historial

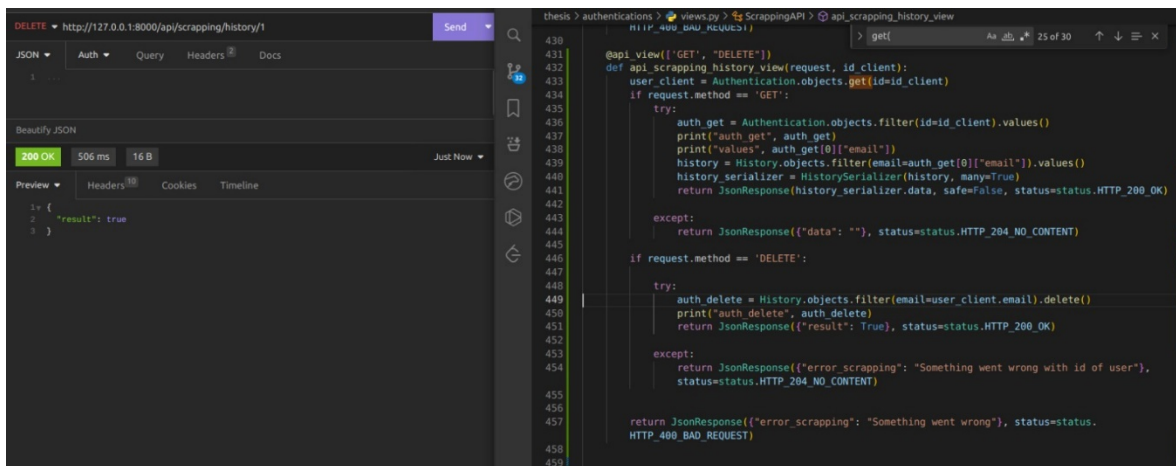
Sprint 17. Eliminación de la API Rest y base de datos

El sprint 17 se enfoca en probar la eliminación de los componentes tanto la API Rest y verificamos la base de datos para probar su funcionamiento. Los resultados de este sprint son:

- Eliminación de la API Rest.

Eliminación de la API Rest

Se procede a realizar la eliminación de la API Rest y luego se verifica el funcionamiento de la base de datos como se muestra en las Fig 43, Fig 44 y Fig 45.



The screenshot shows a web browser on the left and a code editor on the right. The browser displays a DELETE request to `http://127.0.0.1:8000/api/scrapping/history/1` with a 200 OK status. The code editor shows the implementation of the `api_scrapping_history_view` function, which handles GET and DELETE requests. The DELETE method uses `History.objects.filter(email=user_client.email).delete()` to remove the record.

Fig 43. Eliminando la API Rest

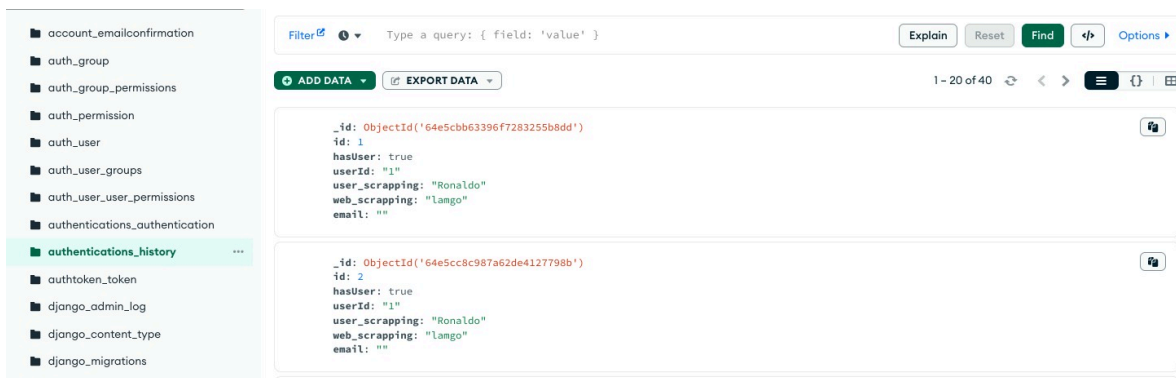


Fig 44. Base de datos antes de eliminar



Fig 45. Base de datos después de eliminar

4 CONCLUSIONES

- Se ha encontrado algunas limitaciones al utilizar la biblioteca instaloader ya que al realizar muchas peticiones la biblioteca hacer redireccionamientos, una manera de no saturarse es limitándose al margen de cuarenta peticiones.
- El funcionamiento del proyecto de acuerdo con los requerimientos establecidos ha sido exitoso, cumpliendo con las funciones principales tanto como el registro y login de los usuarios y le herramienta para realizar el web scraping en las diferentes páginas web.
- La metodología utilizada que en este caso fue RAD (Rapid application developer) fue indispensable para el desarrollo ya que al realizar las pruebas de manera constante según lo que se fue desarrollando permitió que el resultado final sea un software de calidad que solvente todos los requerimientos de una manera práctica y muy rápida.
- La utilización de un framework como Django fue acertada ya que permitió que la app funcione de la mejor manera ya que todos los scripts se realizaron con el lenguaje de programación Python.
- El lenguaje de programación Python apoyó mucho al desarrollo rápido del sistema web ya que con este lenguaje es más fácil realizar webscraping debido a su practicidad y librerías que permiten agilizar mucho más este proceso de extracción de datos.

5 RECOMENDACIONES

- Antes de empezar con el desarrollo es importante revisar las limitaciones que pueden tener algunas librerías y comprobar si los requerimientos necesitan de muchas solicitudes de usuario
- Es recomendable utilizar la metodología RAD para aplicaciones que no requieran enfocarse tanto en la planificación y más en la calidad del software ya que con las pruebas rápidas constantes que se van realizando conforme se van desarrollando cada una de las funcionalidades permite que se ahorre tiempo y los errores se corrijan apenas van apareciendo.
- Es importante saber que el hacer web scraping en algunas ocasiones puede resultar ilegal por las políticas de seguridad que manejan muchos creadores de contenido que se encuentra en la web, por eso es muy importante saber qué tipo de datos se pueden extraer y si estos mismo tengan alguna seguridad anti scrapeo.
- Cuando se realiza web scraping es importante manipular los datos de manera ordenada y estructurada porque de lo contrario la información obtenida no será de mucha ayuda complicando el entendimiento de esta.

6 REFERENCIAS BIBLIOGRÁFICAS

[1 K. T. I. d. reservados, «¿Qué Es el Web Scraping? Cómo Extraer Legalmente el Contenido de la Web,» [En línea]. Available: <https://kinsta.com/es/base-de-conocimiento/que-es-web-scraping/>. [Último acceso: 18 Agosto 2023].

[2 K. I. t. I. d. reservados, «¿Qué Es el Web Scraping? Cómo Extraer Legalmente el Contenido de la Web,» [En línea]. Available: <https://kinsta.com/es/base-de-conocimiento/que-es-web-scraping/>.

[3 M. Martí, «¿Qué es el Web scraping? Introducción y herramientas,» [En línea]. Available: <https://sitelabs.es/web-scraping-introduccion-y-herramientas/>. [Último acceso: 18 Agosto 2023].

[4 S. (. a. L. C. PLANETA DEAGOSTINI FORMACIÓ I UNIVERSITATS ANDORRA, «Metodologías de desarrollo de software,» [En línea]. Available: <https://www.universitatcarlemany.com/actualidad/blog/metodologias-de-desarrollo-de-software/>. [Último acceso: 18 Agosto 2023].

[5 I. Cloud, «¿Qué es el web scraping?,» [En línea]. Available: <https://www.ionos.es/digitalguide/paginas-web/desarrollo-web/que-es-el-web-scraping/>. [Último acceso: 18 Agosto 2023].

[6 T. I. d. r. Crea System, «¿Qué Es Un Sistema Web?,» [En línea]. Available: <https://www.creasystem.net/posts/que-es-un-sistema-web>. [Último acceso: 18 Agosto 2023].

[7 T. I. d. reservados, «Santander Becas,» [En línea]. Available: <https://www.becas-santander.com/es/blog/metodologias-desarrollo-software.html>. [Último acceso: 18 Agosto 2023].

[8 T. I. d. r. Kissflow, «What is Rapid Application Development (RAD)?,» [En línea]. Available: [https://kissflow.com/application-development/rad/rapid-application-development/#:~:text=17.08.2023-,What%20is%20Rapid%20Application%20Development%20\(RAD\)%3F,a%20prototype%2C%20rather%20than%20planning..](https://kissflow.com/application-development/rad/rapid-application-development/#:~:text=17.08.2023-,What%20is%20Rapid%20Application%20Development%20(RAD)%3F,a%20prototype%2C%20rather%20than%20planning..) [Último acceso: 18 Agosto 2023].

[9 K. Beck, «Extreme Programming Explained,» Addison Wesley, 2000, pp. 3-7.
]

[1 V. Solutions, «Qué es la recopilación de requisitos: definición y herramientas,» [En línea].
0] Available: <https://visuresolutions.com/es/blog/requirements-gathering/>. [Último acceso: 19 Agosto 2023].

[1 O. S. t. I. d. reservados, «Arquitectura de software,» [En línea]. Available:
1] <https://openwebinars.net/blog/arquitectura-de-software-que-es-y-que-tipos-existen/>.
[Último acceso: 19 Agosto 2023].

- [1 D. S. Foundation, «Django,» [En línea]. Available: <https://www.djangoproject.com/>. [Último acceso: 19 Agosto 2023].
- [1 T. N. T. School., «¿Qué es Django y para qué se utiliza?,» [En línea]. Available: <https://www.tokioschool.com/noticias/que-es-django/>. [Último acceso: 18 Agosto 2023].
- [1 C. S.L., «Qué es MVC,» [En línea]. Available: [https://desarrolloweb.com/articulos/que-es-4\) mvc.html](https://desarrolloweb.com/articulos/que-es-4) mvc.html). [Último acceso: 18 Agosto 2023].
- [1 R. D. Hernandez, «El patrón modelo-vista-controlador: Arquitectura y frameworks explicados,» [En línea]. Available: <https://www.freecodecamp.org/espanol/news/el-modelo-de-arquitectura-view-controller-pattern/>. [Último acceso: 18 Agosto 2023].
- [1 I. o. s. f. T. I. d. r. Amazon Web Services, «¿Qué es Python?,» [En línea]. Available: <https://aws.amazon.com/es/what-is/python/>. [Último acceso: 18 Agosto 2023].
- [1 MongoDB, «¿Qué es MongoDB?,» [En línea]. Available: [https://www.mongodb.com/es/what-7\) is-mongodb](https://www.mongodb.com/es/what-7) is-mongodb). [Último acceso: 19 Agosto 2023].
- [1 L. Richardson., «Beautiful Soup,» [En línea]. Available: <https://www.crummy.com/software/BeautifulSoup/>. [Último acceso: 19 Agosto 2023].
- [1 A. Caballero, «Descargar Imágenes y Videos desde un Perfil Público de Instagram utilizando 9) Instaloader,» [En línea]. Available: [Descargar_Imagenes_y_Videos_desde_un_Perfil_Publico_de_Instagram_utilizando_Instaloader](#). [Último acceso: 15 Agosto 2023].
- [2 Y. Fernández, «Qué es Github y qué es lo que le ofrece a los desarrolladores,» xataka, 30 0) Octubre 2019. [En línea]. Available: <https://www.xataka.com/basics/que-github-que-que-le-ofrece-a-desarrolladores>. [Último acceso: 24 Mayo 2022].
- [2 Platzi, «Curso de Automatización de Test con Playwright,» [En línea]. Available: <https://platzi.com/cursos/playwright/#:~:text=Playwright%20es%20una%20herramienta%20para,cualquier%20lenguaje%20de%20programaci%C3%B3n%20moderno..> [Último acceso: 20 Agosto 2023].
- [2 «Create API Documentation with Postman,» API PLATFORM, [En línea]. Available: <https://learning.postman.com/docs/getting-started/introduction/>. [Último acceso: 24 Mayo 2022].

7 ANEXOS

ANEXO I. Certificado turnitin

ANEXO I. Manual técnico

ANEXO II. Manual de usuario

ANEXO III. Manual de instalación



**ESCUELA POLITÉCNICA NACIONAL
ESCUELA DE FORMACIÓN DE TECNÓLOGOS
CAMPUS POLITÉCNICO "ING. JOSÉ RUBÉN ORELLANA"**

CERTIFICADO DE ORIGINALIDAD

Quito, D.M. 28 de agosto de 2023

De mi consideración:

Yo, Juan Pablo Zaldumbide Proaño, en calidad de Director del Trabajo de Integración Curricular titulado APLICATIVO WEB DE EXTRACCIÓN DE DATOS asociado al DESARROLLO DE APLICACIÓN WEB PARA EXTRACCIÓN DE DATOS CON WEB SCRAPING Y HERRAMIENTA DE ANALÍTICA DE DATOS elaborado por el estudiante SOLEDISPA MENDOZA CHRISTIAN GABRIEL de la carrera en TECNOLOGÍA SUPERIOR EN DESARROLLO DE SOFTWARE , certifico que he empleado la herramienta Turnitin para la revisión de originalidad del documento escrito completo, producto del Trabajo de Integración Curricular indicado.

El documento escrito tiene un índice de similitud del 10%.

Es todo cuanto puedo certificar en honor a la verdad, pudiendo el interesado hacer uso del presente documento para los trámites de titulación.

NOTA: Se adjunta el informe generado por la herramienta Turnitin.

Atentamente,



firmado electrónicamente por:
**JUAN PABLO
ZALDUMBIDE PROANO**

**Juan Pablo Zaldumbide Proaño
Docente Ocasional a Tiempo Completo
ESFOT**

ANEXO II Manual Técnico

A continuación, se describen con mayor entendimiento los componentes faltantes en el informe técnico del proyecto

1. Recopilación de requerimientos

TABLA II. Recopilación de requerimientos

RECOPIACION DE REQUERMIENTOS		
TIPO DE SISTEMA	ID - RR	ENUNCIADO DEL ITEM
Web para la extracción de datos con web scraping	RR001	Como usuario necesito iniciar sesión.
	RR002	Como usuario necesito acceder al sistema web.
	RR003	Como usuario necesito actualizar mi perfil en la aplicación.
	RR004	Como usuario necesito poder registrarme si no tengo cuenta registrada todavía.
	RR005	Como usuario necesito poder ingresar la página web de la que quiero información.
	RR006	Como usuario necesito ver información sobre lo solicitado.
	RR007	Como usuario necesito poder acceder al a base de datos donde se almacena la información.

2. HISTORIAS DE USUARIO

En las tablas consiguientes de la **TABLA I** hasta la Error! Reference source not found. se muestran la lista de historias de usuario completas del proyecto.

TABLA III Historia de Usuario HU001

HISTORIA DE USUARIO	
Identificador (ID): HU001	Usuario: Administrador, cliente.

Nombre Historia: Iniciar sesión	
Prioridad en Negocio: Alta	Riesgo en Desarrollo: Alta
Iteración Asignada: 3	
Responsable: Christian Soledispa	
Descripción: El administrador y cliente deberán permitir la funcionalidad de inicio de sesión en el <i>sistema web</i> , para que los usuarios puedan iniciar sesión en el sistema.	
Observación: El inicio de sesión deberá contener los siguientes campos: email, contraseña.	

TABLA IV Historia de Usuario HU002

HISTORIA DE USUARIO	
Identificador (ID): HU002	Usuario: Administrador.
Nombre Historia: Mostrar usuarios	
Prioridad en Negocio: Media	Riesgo en Desarrollo: Alta
Iteración Asignada: 2	
Responsable: Christian Soledispa	
Descripción: El administrador deberá tener la funcionalidad de visualizar a los potenciales usuarios finales desde el <i>backend</i> .	
Observación: Los potenciales usuarios deben mostrarse por sus campos: nombre, email.	

TABLA V Historia de Usuario HU003

HISTORIA DE USUARIO	
Identificador (ID): HU003	Usuario: Cliente

Nombre Historia: Editar perfil	
Prioridad en Negocio: Alta	Riesgo en Desarrollo: Media
Iteración Asignada: 3	
Responsable: Christian Soledispa	
Descripción: El usuario como cliente debe tener la opción de actualizar su perfil.	
Observación: Las actualizaciones del perfil del usuario deben ser realizadas en cualquier momento y sin restricción por los mismos dueños del perfil.	

TABLA VI Historia de Usuario HU004

HISTORIA DE USUARIO	
Identificador (ID): HU004	Usuario: Administrador
Nombre Historia: Eliminar usuarios	
Prioridad en Negocio: Media	Riesgo en Desarrollo: Alta
Iteración Asignada: 2	
Responsable: Christian Soledispa	
Descripción: El administrador deberá tener la funcionalidad de eliminar a los usuarios registrados en el <i>backend</i> , dependiendo de cual sea el caso y la razón de ser eliminado.	
Observación: El administrador debe contar con el privilegio de poder eliminar a los usuarios mal intencionados que quieran dañar el sistema.	

TABLA II Historia de Usuario HU005

HISTORIA DE USUARIO	
Identificador (ID): HU005	Usuario: Cliente
Nombre Historia: Elegir la página web	
Prioridad en Negocio: Media	Riesgo en Desarrollo: Alta

Iteración Asignada: 2
Responsable: Christian Soledispa
Descripción: El Cliente podrá elegir antes de ingresar la información a buscar, en cual de ellas quiere realizar la búsqueda.
Observación: Las búsquedas se podrán realizar en las redes sociales más utilizadas como Facebook, Twitter, Tik Tok, Instagram.

TABLA III Historia de Usuario HU006

HISTORIA DE USUARIO	
Identificador (ID): HU006	Usuario: Administrador.
Nombre Historia: Visualizar la información extraída.	
Prioridad en Negocio: Media	Riesgo en Desarrollo: Media
Iteración Asignada: 2	
Responsable: Christian Soledispa	
Descripción: El administrador podrá visualizar por medio del backend la información solicitada por el cliente.	
Observación: La información se mostrará en una interfaz sencilla ya que la misma será proporcionada y mostrada mediante el frontend.	

TABLA IX Historia de Usuario HU007

HISTORIA DE USUARIO	
Identificador (ID): HU007	Usuario: Administrador y cliente
Nombre Historia: Manipular los datos	
Prioridad en Negocio: Media	Riesgo en Desarrollo: Media
Iteración Asignada: 2	

Responsable: Christian Soledispa
Descripción: El administrador podrá manipular y realizar una estructura de los datos extraídos para que sean proporcionados al cliente como tal.
Observación: Los datos serán almacenados en una base de datos NoSQL.

TABLA X Historia de Usuario HU008

HISTORIA DE USUARIO	
Identificador (ID): HU008	Usuario: Administrador y Cliente
Nombre Historia: Acceder a la base de datos	
Prioridad en Negocio: Alta	Riesgo en Desarrollo: Alta
Iteración Asignada: 2	
Responsable: Christian Soledispa	
Descripción: Tanto el Administrador como el cliente podrán acceder a la base de datos para estructurar los datos y puedan ser visualizados.	
Observación: Los datos podrán ser seleccionados para que únicamente sean visualizados los datos que se necesiten.	

ANEXO III Manual de usuario

El siguiente enlace contiene el manual de usuario del sistema, como video en la plataforma youtube:

<https://youtu.be/csEMXSEwu-w>

ANEXO IV Manual de instalación

En el repositorio de GitHub se presenta el manual de cómo instalar y utilizar el sistema, describe el desarrollo completo del proyecto:

https://github.com/ChristianSoledispa/Proyecto_de_titulacion.git