

Webscraping – DIAI

BOUCHIBA Emine

CHENIK Yassine

RAPPORT – PROJET

WEB SCRAPING

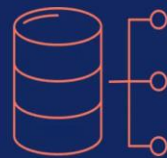


TABLE DES MATIERES

Project Definition	2-3
Webscraping Part.....	4-5
With Excel	4
With BeautifulSoup and Selenium	5
Web Interface	6-9
Getting the startups locations	6-7
Generation of the RSE Scores	8
Streamlit Website Generation	9
Demo Video	10

Our goal :

- ✓ We sincerely believe that there can be no change in society unless we start small, namely with startups. If we change the way we invest, namely by investing in promising startups that are also respectful of RSE standards, it is only then that there can be real effective change. That's why our goal through this project is to build an interface that can provide information about medical startups and, more than that, an interface that can help us choose the startups that are most respectful of RSE.

Project Roadmap :

- ✓ Initially, we planned to scrape data from LinkedIn to gather information about medical startups. However, we soon realized that LinkedIn's data was insufficient for our interface, focusing more on business-related details rather than practical information.
- ✓ Consequently, we shifted our approach to constructing our own list of medical startups. We explored alternative data scraping techniques beyond BeautifulSoup and Selenium, opting to use Excel. This method involved scraping data from various medical websites known for their extensive startup information, such as Biocat, Biopartner, Bionow, Atlantapole, ci3, and others.
- ✓ Having obtained data from these websites, we then employed a combination of BeautifulSoup and Selenium to scrape URL links from our previously created Excel sheet. This enabled us to extract information directly from the startups' websites.
- ✓ With this approach, we successfully compiled a comprehensive dataset containing all the necessary information. Our next focus was on developing the web interface.
- ✓ For the interface, we chose to use Streamlit, diverging from our usual tools to explore new methods. Streamlit proved to be an efficient choice, requiring only a single Python script instead of separate HTML and CSS files. This streamlined the process of creating faster, more customizable websites.

- ✓ With the web interface solution decided, we then considered which features to include :
 - **1st feature** : Displaying information about all startups, filterable by the « sector » column.
 - **2nd feature** : Visualizing startup locations using the « geopy » library to obtain « latitude » and « longitude » data from string inputs.
 - **3rd feature** : Presenting each startup's RSE score, filtered by the 'sector' column. This involved using NLP models from Hugging Face to identify words in the « specialization » column that closely match the RSE lexicon.

Webscraping Part

With Excel :

✓ We used this function in order to scrap the data :

The screenshot shows an Excel spreadsheet with the following data:

Name	Link	Address line1	Address line2	Address line3	Phone number	Summary
A&O Pharmad	www.aopharm.de	Am Sattel 17	79588 Efringen-Kirch	07628 95 03 1...	...	Diese E-Ma Arzneimittel freigegeben klinische Prüfpräparate und Marktware durch die Sachkundige Person mit eig...
Albert-Ludwig	www.informa.de	Institut für Inf...	79110 Freiburg	0761 203-746...	...	Diese E-Ma Arcondis ist eine Unternehmensberatung für das Management von Information, Qualität und der IT...
Arcondis AG	www.arcondis.de	Christoph Mer	0 Reinach	0041 61 717 8...	...	Diese E-Ma BESTMINDS, die Personalberatung mit Hauptsitz in Freiburg, hat sich auf die Branchen Life Scienc...
ATG biosynth	www.atg-bios.de	Weberstr. 40	79249 Merzhausen	0761 888 94 2...	...	Diese E-Ma Ich suche Fach- und Führungskräfte aus allen Fachbereichen und Disziplinen. Mehr als 20 Jahre E...
BBI Solutions	www.biarect.de	Bötzingen Str.	79111 Freiburg	0761 47979 0...	...	Diese E-Ma Wir sind ein kleines unabhä... Zu unseren Kunden gehören Qualität, basierend auf unserem langjahr...
BESTMINDS	www.bestminds.de	Basler Str. 65	79100 Freiburg	+49-761-888 5...	...	Diese E-Ma Die BioCopy AG ist ein junges Biotech-Startup mit Hauptsitz bei Basel, Schweiz und einer Forschun...
BiochemA Gm	www.biochem.de	Im Oberwald 1	79359 Riegel	07642 7018...	...	Diese E-Ma BioCopy's vielfach ausgezeichnetes Team von mehr als 20 Experten aus den Bereichen Biologie, P...
BioCopy Gmb	www.biocopy.de	Elzstr. 27	79312 Emmendingen-	Diese E-Ma Die BioFluidix GmbH hat sich als Unternehmen aus dem Bereich der Mikrosystemtechnik auf die Ha...
BioFluidix Gm	www.biofluidix.de	Engesserstr. 4	79108 Freiburg	0761 458 938...	...	Diese E-Ma Als innovatives Unternehmen haben wir uns der kontinuierlichen Weiterentwicklung unserer Produk...
BIOSS Centre	www.bioss.uni-freiburg.de	Schänzlestr. 1	79104 Freiburg	0761 203 97 3...	...	Diese E-Ma Biologische Si Lösung wichtig und ein weser Fortschritt der... Das Centre for Biological Signalling Stu...
CapCo Bio Gr	www.capcobio.de	Egonstr. 51-5	79106 Freiburg	0761 8884220...	...	Diese E-Ma Die CapCo Bi-SNKS... Mit Hilfe dieser biokompatiblen Kapseln können Biomoleküle in leber...
CellGenix Gm	www.cellgenix.de	Am Flughafen	79108 Freiburg	0761 88889 21...	...	Diese E-Ma Der Schwerpunkt von CellGenix ist die Entwicklung und GMP-Herstellung hochwertiger Reagenzien...
Charles River	www.crivier.com	Am Flughafen	79108 Freiburg	0761 51559 1f...	...	Diese E-Ma Die Entwicklung und Herstellung der Produkte findet in modernsten GMP-Reinraumlaboren in Freib...
ChemCon Gm	www.chemcon.de	Engesserstr. 4	79108 Freiburg	0761 5597 44f...	...	Diese E-Ma Charles River in vitro Assay-System in vivo Wirksamkeits in vivo Arbeiten sir...
Chemingenier	www.chemingenier.de	Binnigerstr. 2	0 Münchenstein	0041 61 467 5...	...	Diese E-Ma Die ChemCon GmbH ist auf die organische und anorganische Chemie kleiner Moleküle spezialisiert...
CRC Clean R	www.crc.de	Badenweilerst	79115 Freiburg	0761 4 78 13 f...	...	Diese E-Ma Chemingenier bietet hochwertige und praxiserprobte und Ingenieur- und Beratungsdienstleistungen...
CW-NOTIO D	www.cw-notio.de	Jägerhauslewi	79104 Freiburg	0761 208 92 7...	...	Diese E-Ma Die CRC Clean Room Consulting ist ein international tätiges Planungsbüro für High-Tech Industrie...
Cytiva Europe	www.cytiva.com	Munzinger Str	79111 Freiburg	0761 600 49 7...	...	Diese E-Ma Unsere Leistungen im Überblick:
Deutsche Ban	www.db.com	Rotteckring 3	79098 Freiburg	0761 2184 24...	...	Diese E-Ma Wir bei Cytiva glauben, dass Erfahrung Sie mehr unter cytiva.com.
Drees & Somr	www.drees-sommer.de	Salzstraße 15	79098 Freiburg	0711 1317 20...	...	Diese E-Ma Ob Neubau oder Revitalisierungen: Drees & Sommer verbindet jahrzehntelange Erfahrung bei der F...
DSM Nutrition	www.dsm.com	Emil-Barell-Str	79630 Grenzach-Wyl-	Diese E-Ma Mit individuellen Konzepten und professioneller Umsetzung sichern wir Ihnen Qualität, Kosten und T...
HE	www.he.com	Unterbaselwe	79576 Weil am Rhei	07621 798373...	...	Diese E-Ma Wir bieten ein Leistungsspektrum vom Controlling bis hin zum Rundum-Sorglos-Paket mit Kosten-...

Webscrapping Part

With BeautifulSoup and Selenium :

✓ The code is in the notebook, however this is the given .xlsx file :

Enregistrement automatique

Startup database... • Enregistré dans ce PC

Rechercher

Fichier

Accueil

Insertion

Dessin

Mise en page

Formules

Données

Révision

Affichage

Automate

Aide

Commentaires

Partager

Coller

Presse-papiers

Calibri

11

A

A

G

I

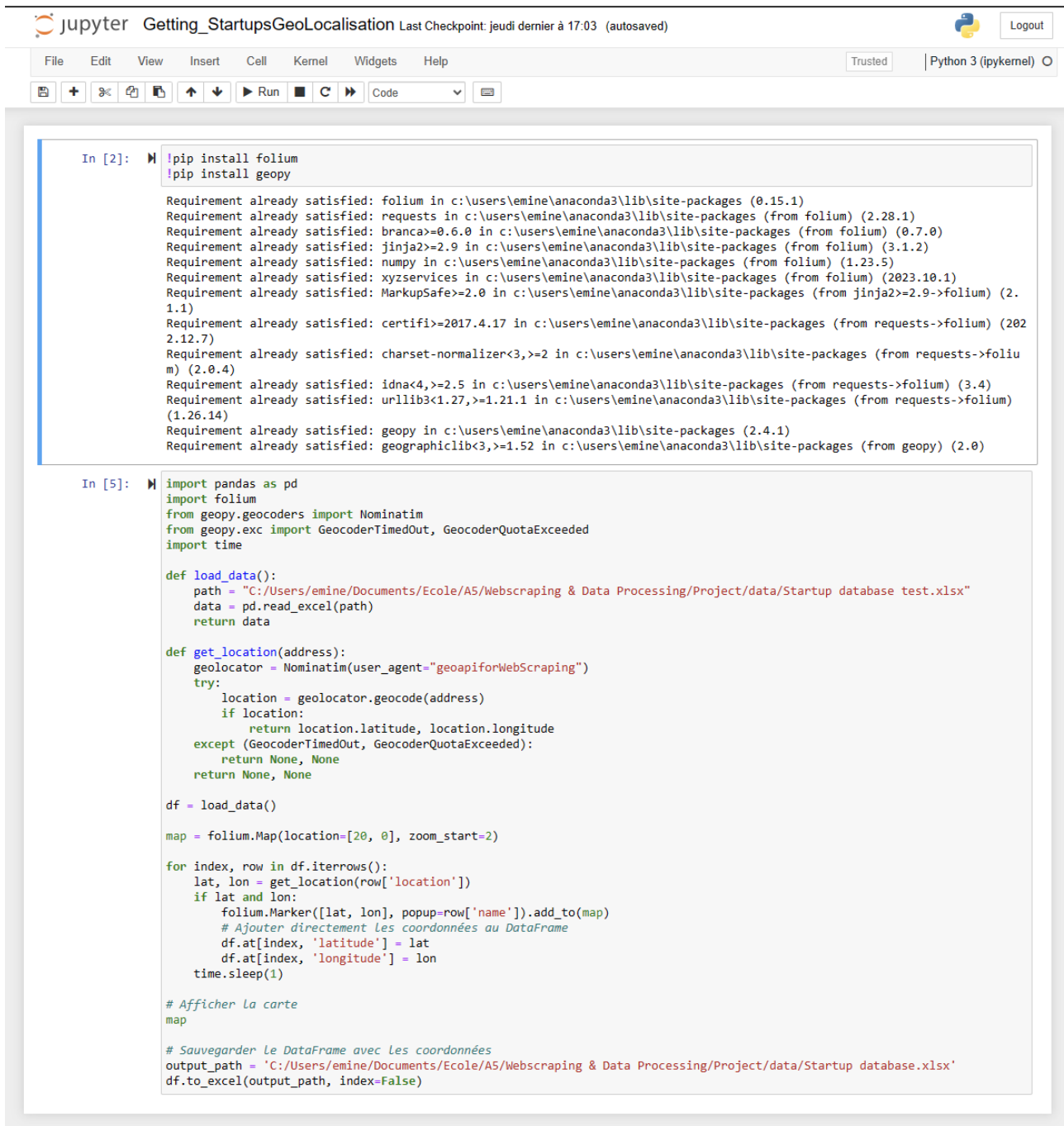
S

</

Web Interface

Getting the startups locations :

- ✓ Using the GeoPy Library in order to get the startups exact location from an str input :



Jupyter Getting_StartupsGeoLocalisation Last Checkpoint: jeudi dernier à 17:03 (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (ipykernel)

```

In [2]: !pip install folium
        !pip install geopy

Requirement already satisfied: folium in c:\users\emine\anaconda3\lib\site-packages (0.15.1)
Requirement already satisfied: requests in c:\users\emine\anaconda3\lib\site-packages (from folium) (2.28.1)
Requirement already satisfied: branca>=0.6.0 in c:\users\emine\anaconda3\lib\site-packages (from folium) (0.7.0)
Requirement already satisfied: Jinja2>=2.9 in c:\users\emine\anaconda3\lib\site-packages (from folium) (3.1.2)
Requirement already satisfied: numpy in c:\users\emine\anaconda3\lib\site-packages (from folium) (1.23.5)
Requirement already satisfied: xyzservices in c:\users\emine\anaconda3\lib\site-packages (from folium) (2023.10.1)
Requirement already satisfied: MarkupSafe>=2.0 in c:\users\emine\anaconda3\lib\site-packages (from Jinja2>=2.9->folium) (2.1.1)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\emine\anaconda3\lib\site-packages (from requests->folium) (2022.12.7)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\users\emine\anaconda3\lib\site-packages (from requests->folium) (2.0.4)
Requirement already satisfied: idna<4,>=2.5 in c:\users\emine\anaconda3\lib\site-packages (from requests->folium) (3.4)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\users\emine\anaconda3\lib\site-packages (from requests->folium) (1.26.14)
Requirement already satisfied: geopy in c:\users\emine\anaconda3\lib\site-packages (2.4.1)
Requirement already satisfied: geographiclib<3,>=1.52 in c:\users\emine\anaconda3\lib\site-packages (from geopy) (2.0)

In [5]: import pandas as pd
import folium
from geopy.geocoders import Nominatim
from geopy.exc import GeocoderTimedOut, GeocoderQuotaExceeded
import time

def load_data():
    path = "C:/Users/emine/Documents/Ecole/A5/Webscraping & Data Processing/Project/data/Startup database test.xlsx"
    data = pd.read_excel(path)
    return data

def get_location(address):
    geolocator = Nominatim(user_agent="geopapirforWebScraping")
    try:
        location = geolocator.geocode(address)
        if location:
            return location.latitude, location.longitude
    except (GeocoderTimedOut, GeocoderQuotaExceeded):
        return None, None
    return None, None

df = load_data()

map = folium.Map(location=[20, 0], zoom_start=2)

for index, row in df.iterrows():
    lat, lon = get_location(row['location'])
    if lat and lon:
        folium.Marker([lat, lon], popup=row['name']).add_to(map)
        # Ajouter directement Les coordonnées au DataFrame
        df.at[index, 'latitude'] = lat
        df.at[index, 'longitude'] = lon
    time.sleep(1)

# Afficher La carte
map

# Sauvegarder Le DataFrame avec Les coordonnées
output_path = 'C:/Users/emine/Documents/Ecole/A5/Webscraping & Data Processing/Project/data/Startup database.xlsx'
df.to_excel(output_path, index=False)
  
```


Web Interface

Displaying the startups locations :

✓ Using the folium library from python :

jupyter Displaying_StartupsGeoLocalisation Last Checkpoint: jeudi dernier à 18:00 (autosaved) Python 3 (pykernel)

File Edit View Insert Cell Kernel Widgets Help Trusted | Python 3 (pykernel)

```
In [1]: import pandas as pd
import folium

# Fonction pour charger Les données depuis Le fichier Excel
def load_coordinates():
    path = "C:\\Users\\emine\\Documents\\Ecole\\AS\\Webcraping & Data Processing\\Project\\data\\startup_coordinates.xlsx"
    data = pd.read_excel(path)
    return data


# Charger Les données de coordonnées
df_coordinates = load_coordinates()

# Créer La carte
map = folium.Map(location=[20, 0], zoom_start=2)

# Ajouter des marqueurs pour chaque startup
for _, row in df_coordinates.iterrows():
    folium.Marker([row['latitude'], row['longitude']], popup=row['name']).add_to(map)

# Afficher La carte avec Streamlit
map
```

Out[1]:



```
In [ ]: import pandas as pd
import folium

def load_data():
    path = "C:\\Users\\emine\\Documents\\Ecole\\AS\\Webcraping & Data Processing\\Project\\data\\Startup database final.xlsx"
    data = pd.read_excel(path)
    return data

# Obtenir Les coordonnées géographiques
def get_location(address):
    # Ici, vous pouvez ajouter votre logique de géocodage
    # ou utiliser une méthode pour extraire Les coordonnées si elles sont déjà présentes dans Les données
    pass

# Appel de La fonction de chargement de données
df = load_data()

# Sélectionner un secteur via un menu déroulant
sector = st.selectbox('Choisissez un secteur', df['sector'].unique())

# Filtrer Les données par secteur
filtered_data = df[df['sector'] == sector]

# Afficher Les données filtrées
st.subheader(f'Données pour le secteur : {sector}')
st.write(filtered_data[['name', 'phone', 'size', 'website', 'founded', 'location']])

# Création de La carte pour Le secteur sélectionné
map = folium.Map(location=[20, 0], zoom_start=2)
for _, row in filtered_data.iterrows():
    lat, lon = get_location(row['location'])
    if lat and lon:
        folium.Marker([lat, lon], popup=row['name']).add_to(map)

# Fonction pour charger Les données depuis Le fichier Excel
def load_coordinates():
    path = "C:\\Users\\emine\\Documents\\Ecole\\AS\\Webcraping & Data Processing\\Project\\data\\startup_coordinates.xlsx"
    data = pd.read_excel(path)
    return data

# Charger Les données de coordonnées
df_coordinates = load_coordinates()

# Créer La carte
map = folium.Map(location=[20, 0], zoom_start=2)

# Ajouter des marqueurs pour chaque startup
for _, row in df_coordinates.iterrows():
    folium.Marker([row['latitude'], row['longitude']], popup=row['name']).add_to(map)
```


Web Interface

Generation of the RSE Scores :

- ✓ Using the transformers library created by Hugging Face for NLP integration in python :

jupyter RSE_Score Last Checkpoint: il y a 15 heures (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel)

```

In [1]: !pip install sentence_transformers
import pandas as pd

# Charger les données
def load_data():
    path = "C:/Users/emine/Documents/Ecole/A5/Webscraping & Data Processing/Project/data/Startup database.xlsx"
    return pd.read_excel(path)

df = load_data()

# Extraire les mots de la colonne 'specialization'
all_words = set()
for specialization in df['specialisation'].dropna():
    words = specialization.split()
    all_words.update(words)

# Convertir l'ensemble en liste pour faciliter l'utilisation
unique_words = list(all_words)

# Afficher les mots uniques
print(unique_words)

```

Requirement already satisfied: sentence_transformers in c:\users\emine\anaconda3\lib\site-packages (2.2.2)
Requirement already satisfied: torch>=1.6.0 in c:\users\emine\anaconda3\lib\site-packages (from sentence_transformers) (2.1.2)
Requirement already satisfied: sentencepiece in c:\users\emine\anaconda3\lib\site-packages (from sentence_transformers) (0.1.99)
Requirement already satisfied: huggingface-hub>=0.4.0 in c:\users\emine\anaconda3\lib\site-packages (from sentence_transformers) (0.10.1)
Requirement already satisfied: scipy in c:\users\emine\anaconda3\lib\site-packages (from sentence_transformers) (1.10.0)
Requirement already satisfied: scikit-learn in c:\users\emine\anaconda3\lib\site-packages (from sentence_transformers) (1.2.1)
Requirement already satisfied: torchvision in c:\users\emine\anaconda3\lib\site-packages (from sentence_transformers) (0.16.2)
Requirement already satisfied: transformers<5.0.0,>=4.6.0 in c:\users\emine\anaconda3\lib\site-packages (from sentence_transformers) (4.24.0)
Requirement already satisfied: nltk in c:\users\emine\anaconda3\lib\site-packages (from sentence_transformers) (3.7)
Requirement already satisfied: tqdm in c:\users\emine\anaconda3\lib\site-packages (from sentence_transformers) (4.64.1)
Requirement already satisfied: numpy in c:\users\emine\anaconda3\lib\site-packages (from sentence_transformers) (1.23.5)
Requirement already satisfied: typing-extensions>=3.7.4.3 in c:\users\emine\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence_transformers) (4.9.0)

```

In [2]: from transformers import AutoTokenizer, AutoModelForSequenceClassification
from scipy.spatial.distance import cosine
from sentence_transformers import SentenceTransformer
import numpy as np

# Initialisation du modèle de similarité sémantique
model = SentenceTransformer('sentence-transformers/all-MiniLM-L6-v2')

# Liste de mots-clés RSE
rse_keywords = ["sustainability", "ethical", "social", "environment", "health", "well-being", "community",
               "green", "renewable", "eco-friendly", "inclusive", "diversity", "equality", "charity", "volunteer"]

# Convertir les listes de mots en embeddings
rse_embeddings = model.encode(rse_keywords)
input_embeddings = model.encode(unique_words)

# Trouver les mots les plus similaires
threshold = 0.65 # Seuil de similarité, ajustable selon les besoins
rse_related_words = []

for word, word_embedding in zip(unique_words, input_embeddings):
    similarities = [1 - cosine(word_embedding, rse_embedding) for rse_embedding in rse_embeddings]
    if max(similarities) > threshold:
        rse_related_words.append(word)

print(rse_related_words)

```

['Wellness', 'wellbeing', 'Health', 'Well-being', 'health', 'Wellbeing', 'Community', 'Medical', 'healthcare', 'Communities', 'wellness', 'Health', 'Environmental', 'sustainability', 'Medicine', 'healthcare', 'sustainable', 'Wellness', 'wellness', 'medical', 'Healthcare', 'medicine', 'Fundraising', 'wellbeing', 'Healthcare', 'Communities', 'health']

Web Interface

Streamlit Website Generation :

- ✓ Installation of the necessary libraries :

```
pip install streamlit pandas openpyxl matplotlib seaborn folium streamlit-folium geopy streamlit_option_menu scipy sentence_transformers
```

- ✓ Running Streamlit locally :

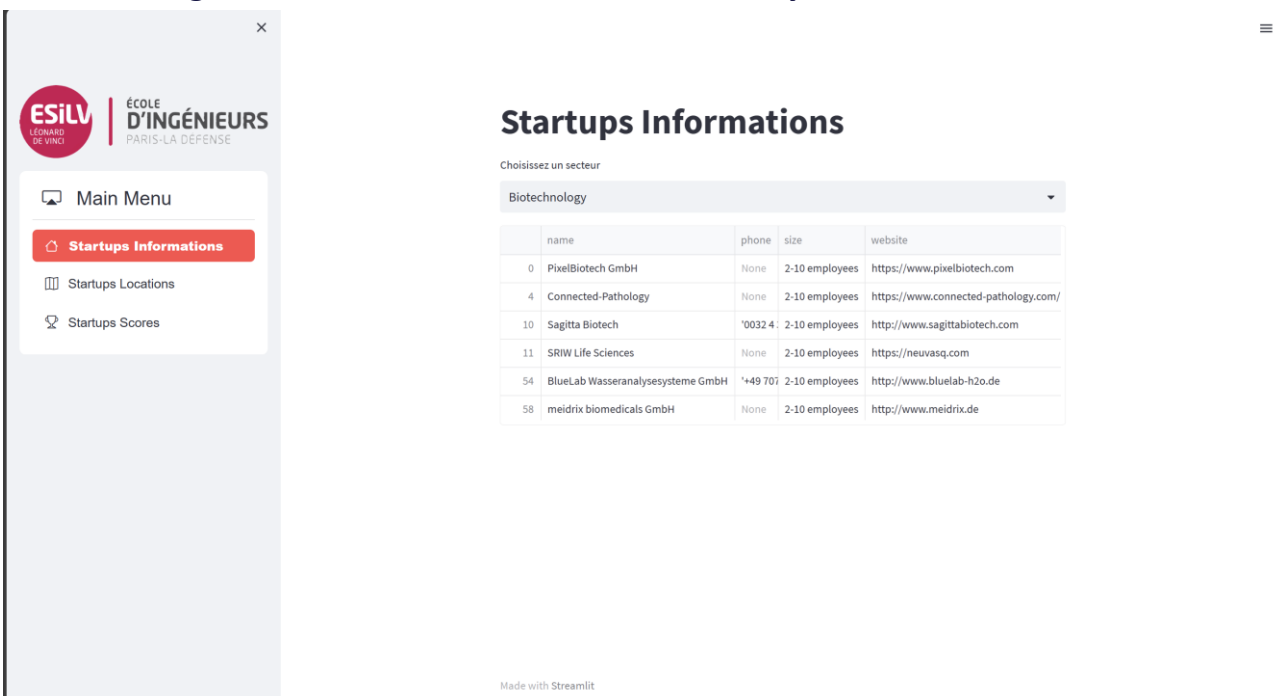
```
Microsoft Windows [version 10.0.22621.2861]
(c) Microsoft Corporation. Tous droits réservés.

C:\Users\emine>streamlit run "C:\Users\emine\Documents\Ecole\A5\Web scraping & Data Processing\Project\code\Project_WebInterface.py"

You can now view your Streamlit app in your browser.

Local URL: http://localhost:8501
Network URL: http://192.168.1.26:8501
```

- ✓ Having the web interface launched and ready to be used :



Streamlit Website Demo :

✓ [Demo Video](#)