

# Leboncoin Scraper

Aimad Hamdaoui

22 Décembre 2024

## Table des matières

<b>1</b>	<b>Présentation générale</b>	<b>1</b>
<b>2</b>	<b>Pré-requis</b>	<b>1</b>
<b>3</b>	<b>Installation</b>	<b>2</b>
<b>4</b>	<b>Configuration des fichiers</b>	<b>2</b>
4.1	Fichier <code>keywords.txt</code> . . . . .	2
4.2	Fichier <code>database.json</code> . . . . .	2
<b>5</b>	<b>Fonctionnement du script</b>	<b>3</b>
5.1	Lecture des mots-clés . . . . .	6
5.2	Chargement de la base de données locale . . . . .	6
5.3	Lancement du navigateur avec <code>undetected-playwright</code> . . . . .	6
5.4	Collecte et traitement des annonces . . . . .	7
5.5	Filtrage des annonces par mots-clés . . . . .	7
5.6	Envoi vers Discord . . . . .	7
5.7	Mise à jour de la base de données . . . . .	7
<b>6</b>	<b>Exécution du script</b>	<b>7</b>
<b>7</b>	<b>Conseils et conclusions</b>	<b>7</b>

## 1 Présentation générale

Le **Leboncoin Scraper** est un script en Python conçu pour automatiser la récupération d'annonces sur le site Leboncoin. Il utilise la bibliothèque `undetected-playwright` pour exécuter un navigateur Chromium, contourner d'éventuelles mesures anti-bots et extraire les informations d'annonces automobiles. Il envoie ensuite les nouvelles annonces correspondant à des mots-clés définis vers un **webhook Discord**.

## 2 Pré-requis

Avant de commencer, assurez-vous de disposer des éléments suivants :

- **Python 3.7+** installé sur votre machine ;
- Un **webhook Discord** valide (la configuration du webhook se fait dans les Paramètres de votre serveur Discord, dans la section *Intégrations → Webhooks*) ;
- La capacité d’installer et d’exécuter **Playwright** (en l’occurrence `undetected-playwright`) ;
- Un accès internet stable.

## 3 Installation

1. **Cloner ou télécharger** ce projet depuis le dépôt Git ou le dossier qui contient le script `.py`.

2. **Créer un environnement virtuel** (optionnel, mais recommandé) :

```
1 python -m venv venv
2 source venv/bin/activate    # Sur Linux/Mac
3 venv\Scripts\activate.bat  # Sur Windows
```

3. **Installer les dépendances** (notamment `requests`, `undetected-playwright`, etc.) :

```
1 pip install -r requirements.txt
```

4. **Installer les navigateurs** via Playwright :

```
1 playwright install
```

## 4 Configuration des fichiers

### 4.1 Fichier `keywords.txt`

Le fichier `keywords.txt` doit contenir la liste des **mots-clés** que vous souhaitez rechercher sur Leboncoin. Chaque mot-clé doit être sur une ligne distincte. Par exemple :

```
peugeot
renault
clio
golf
```

Le script va lire ce fichier au démarrage et ne prendra en compte que les annonces dont le titre contient au moins l’un de ces mots-clés (en minuscules ou majuscules, car la vérification est insensible à la casse).

### 4.2 Fichier `database.json`

Le fichier `database.json` sert de base de données locale pour stocker les liens des annonces déjà consultées ou déjà envoyées sur Discord. Ainsi, le script évite d’envoyer plusieurs fois la même annonce. S’il n’existe pas, il sera créé automatiquement lors de la première exécution du script.

## 5 Fonctionnement du script

Le script principal est un fichier .py qui utilise les bibliothèques `asyncio`, `undetected-playwright`, `requests` et `json`. Ci-dessous, un exemple de structure du code :

```
1 import asyncio
2 import json
3 import requests
4 from pathlib import Path
5 from undetected_playwright.async_api import async_playwright
6
7 LEBONCOIN_SEARCH_URL = "https://www.leboncoin.fr/recherche?
8     category=2&fuel=2&price=min-3500&mileage=min-250000&sort=time"
9 DISCORD_WEBHOOK_URL = "https://discord.com/api/webhooks/..."
10 DB_FILE = Path("database.json")
11 KEYWORDS_FILE = Path("keywords.txt")
12
13 def lire_keywords():
14     """Lit les mots-clés depuis le fichier keywords.txt."""
15     if not KEYWORDS_FILE.exists():
16         print(f"Le fichier {KEYWORDS_FILE} introuvable. Veuillez le
17             créer et y ajouter des mots-clés, un par ligne.")
18         return []
19     with KEYWORDS_FILE.open("r", encoding="utf-8") as f:
20         return [line.strip().lower() for line in f if line.strip()]
21
22 def envoyer_vers_discord(titre, prix, annee, kilometrage, lien):
23     """Envoie l'annonce sur Discord via un webhook."""
24     try:
25         embed = {
26             "title": titre,
27             "description": (
28                 f"**Prix** : {prix}\n"
29                 f"**Année** : {annee}\n"
30                 f"**Kilométrage** : {kilometrage}\n\n"
31                 f"[Voir l'annonce]({lien})"
32             ),
33             "color": 0x7289DA,
34         }
35         data = {
36             "content": "Nouvelle annonce trouvée !",
37             "embeds": [embed],
38         }
39         response = requests.post(DISCORD_WEBHOOK_URL, json=data)
40         response.raise_for_status()
41         print("Annonce envoyée sur Discord avec succès.")
42     except requests.RequestException as e:
43         print(f"Erreur lors de l'envoi sur Discord : {e}")
44
45 async def main():
46     """Relance la page de recherche toutes les 20 secondes après
47     l'envoi sur Discord.
```

```

    avoir fini de traiter les annonces."""
45 keywords = lire_keywords()
46 if not keywords:
47     print("Aucun mot-clé trouvé. Le script s'arrête.")
48     return
49
50 # Lecture (ou création) de la base de données locale
51 if DB_FILE.exists():
52     with DB_FILE.open("r", encoding="utf-8") as f:
53         db = json.load(f)
54 else:
55     db = []
56
57 async with async_playwright() as p:
58     args = ["--disable-blink-features=AutomationControlled"]
59     browser = await p.chromium.launch(headless=False, args=args)
60     context = await browser.new_context(
61         user_agent=(
62             "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91.0.4472.124 Safari/537.36"
63         ),
64         viewport={"width": 1280, "height": 720}
65     )
66     page = await context.new_page()
67
68     try:
69         while True: # Boucle infinie
70             try:
71                 print("Navigation vers la page de recherche")
72                 ...
73                 await page.goto(LEBONCOIN_SEARCH_URL, timeout=60000)
74                 await page.wait_for_selector("a[data-test-id='ad']", timeout=60000)
75
76                 # Récupération de toutes les annonces
77                 annonces = await page.query_selector_all("a[data-test-id='ad']")
78                 print(f"Nombre total d'annonces trouvées : {len(annonces)}")
79
80                 for annonce in annonces:
81                     try:
82                         # Titre
83                         titre_elem = await annonce.
84                         query_selector("p[data-qa-id='aditem_title']")
85                         if not titre_elem:
86                             continue

```

```

87     titre = await titre_elem.inner_text()
88
89     # V rifie si le titre contient l'un
90     # des mots-cl s
91     if not any(keyword in titre.lower()
92         for keyword in keywords):
93         continue
94
95     # Lien
96     lien = await annonce.get_attribute("
97         href")
98     if not lien:
99         continue
100    lien_complet = f"https://www.
101        leboncoin.fr{lien}"
102
103    # Prix
104    prix_elem = await annonce.
105        query_selector("p[data-test-id='
106            price']")
107    prix = await prix_elem.inner_text()
108    if prix_elem else "Non sp cifi "
109
110    # Ann e et kilom trage
111    divs = await annonce.
112        query_selector_all("div.relative.h-
113            full.whitespace-nowrap")
114    if len(divs) >= 2:
115        # Ann e
116        annee_elem = await divs[0].
117            query_selector("p:last-child")
118        annee = await annee_elem.
119            inner_text() if annee_elem else
120                "Non sp cifi e"
121
122        # Kilom trage
123        km_elem = await divs[1].
124            query_selector("p:last-child")
125        kilometrage = await km_elem.
126            inner_text() if km_elem else "
127                Non sp cifi "
128
129    else:
130        annee = "Non sp cifi e"
131        kilometrage = "Non sp cifi "
132
133    # V rifie si d j pr sent dans la
134    # DB
135    if lien_complet in db:
136        continue
137
138    # Ajout dans la DB et envoi Discord

```

```

122         db.append(lien_complet)
123         envoyer_vers_discord(titre, prix,
124                               annee, kilometrage, lien_complet)
125
126         print(f"Nouvelle annonce ajout e : {titre} | {prix} | {annee} | {kilometrage}")
127     except Exception as e:
128         print(f"Erreur de traitement d'une
129               annonce : {e}")
130
131     # Sauvegarde de la DB
132     with DB_FILE.open("w", encoding="utf-8") as f
133     :
134         json.dump(db, f, ensure_ascii=False,
135                   indent=4)
136
137         print("Termin pour cette boucle. Relance de
138               la recherche dans 20 secondes...")
139         await asyncio.sleep(20)
140
141     except Exception as e:
142         print(f"Erreur lors de la navigation ou du
143               parsing : {e}")
144         await asyncio.sleep(20)
145
146     finally:
147         await browser.close()
148
149 if __name__ == "__main__":
150     asyncio.run(main())

```

Voici les différentes étapes clés :

## 5.1 Lecture des mots-clés

La fonction lirekeywords() ouvre le fichier keywords.txt et en charge les mots-clés, un par ligne. Chaque mot-clé est converti en minuscules pour simplifier la vérification.

## 5.2 Chargement de la base de données locale

Le script vérifie l'existence d'un fichier database.json. S'il existe, il est chargé en mémoire pour récupérer la liste des liens déjà envoyés. S'il n'existe pas, une liste vide est initialisée.

## 5.3 Lancement du navigateur avec undetected-playwright

Le script lance une instance de Chromium en mode headless=False (navigateur visible). Des arguments sont passés pour contourner les techniques de détection de bot. Une page est ensuite ouverte dans ce contexte.

## 5.4 Collecte et traitement des annonces

Le script navigue vers l'URL de recherche Leboncoin (stocker dans `LEBONCOIN_SEARCH_URL`). Une fois la page chargée, il recherche toutes les balises correspondantes aux annonces (`a[data-test-id='ad']`). Pour chacune d'entre elles :

- Il récupère le titre, le lien, le prix, l'année et le kilométrage ;
- Il vérifie si le titre contient l'un des mots-clés.

## 5.5 Filtrage des annonces par mots-clés

Le script ne conserve que les annonces dont le titre contient au moins un mot-clé du fichier `keywords.txt`. La correspondance se fait en minuscules.

## 5.6 Envoi vers Discord

Si l'annonce n'est pas déjà dans la base de données (`database.json`), la fonction `envoyer_vers_discord()` est appelée. Elle construit un objet `embed` (spécifique à Discord) et poste une requête HTTP vers le webhook défini dans `DISCORD_WEBHOOK_URL`.

## 5.7 Mise à jour de la base de données

Le lien de l'annonce est ensuite ajouté à la liste en mémoire, puis la base de données est mise à jour en écrivant dans `database.json`. Cela empêche toute future duplication d'annonces.

# 6 Exécution du script

Pour exécuter le script, il suffit de lancer :

```
1 python votre_script.py
```

Assurez-vous d'avoir bien configuré vos fichiers `keywords.txt` et `database.json` (ou de laisser le script créer le second) ainsi que votre `DISCORD_WEBHOOK_URL` dans le code.

Le script reste en boucle **infinie** ; toutes les 20 secondes, il recharge la page, inspecte les annonces, envoie les nouvelles annonces correspondantes, puis attend de nouveau 20 secondes.

# 7 Conseils et conclusions

- Le script étant dans une boucle infinie, vous pouvez l'arrêter manuellement avec `CTRL + C`.
- Vous pouvez adapter l'URL de recherche `LEBONCOIN_SEARCH_URL` selon vos propres critères (catégorie, filtre de prix, kilométrage, etc.).
- Ajustez la `DISCORD_WEBHOOK_URL` avec celle de votre propre serveur Discord.
- N'hésitez pas à personnaliser la fréquence de mise à jour (actuellement 20 secondes via `asyncio.sleep(20)`).
- Si vous exécutez le navigateur en `headless=True`, vous ne verrez pas l'interface, mais cela fonctionnera en tâche de fond.